

# An Eagle's Eye View of XML

---

This first chapter introduces you to XML. It explains in general what XML is and how it is used. It shows you how the different pieces of the XML equation fit together, and how an XML document is created and delivered to readers.

## What Is XML?

XML stands for Extensible Markup Language (often written as eXtensibleMarkup Language to justify the acronym). XML is a set of rules for defining semantic tags that break a document into parts and identify the different parts of the document. It is a meta-markup language that defines a syntax used to define other domain-specific, semantic, structured markup languages.

## XML Is a Meta-Markup Language

The first thing you need to understand about XML is that it isn't just another markup language like the Hypertext Markup Language (HTML) or troff. These languages define a fixed set of tags that describe a fixed number of elements. If the markup language you use doesn't contain the tag you need — you're out of luck. You can wait for the next version of the markup language hoping that it includes the tag you need; but then you're really at the mercy of what the vendor chooses to include.

XML, however, is a meta-markup language. It's a language in which you make up the tags you need as you go along. These tags must be organized according to certain general principles, but they're quite flexible in their meaning. For instance, if you're working on genealogy and need to describe people, births, deaths, burial sites, families, marriages, divorces, and so on, you can create tags for each of these. You don't have to force your data to fit into paragraphs, list items, strong emphasis, or other very general categories.



### In This Chapter

What is XML?

Why are developers excited about XML?

The life of an XML document

Related technologies



The tags you create can be documented in a Document Type Definition (DTD). You'll learn more about DTDs in Part II of this book. For now, think of a DTD as a vocabulary and a syntax for certain kinds of documents. For example, the MOL.DTD in Peter Murray-Rust's Chemical Markup Language (CML) describes a vocabulary and a syntax for the molecular sciences: chemistry, crystallography, solid state physics, and the like. It includes tags for atoms, molecules, bonds, spectra, and so on. This DTD can be shared by many different people in the molecular sciences field. Other DTDs are available for other fields, and you can also create your own.

XML defines a meta syntax that domain-specific markup languages like MusicML, MathML, and CML must follow. If an application understands this meta syntax, it automatically understands all the languages built from this meta language. A browser does not need to know in advance each and every tag that might be used by thousands of different markup languages. Instead it discovers the tags used by any given document as it reads the document or its DTD. The detailed instructions about how to display the content of these tags are provided in a separate style sheet that is attached to the document.

For example, consider Schrodinger's equation:

$$i\hbar \frac{\partial \psi(\mathbf{r}, t)}{\partial t} = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi(\mathbf{r}, t)}{\partial x^2} + V(\mathbf{r}) \psi(\mathbf{r}, t)$$

Scientific papers are full of equations like this, but scientists have been waiting eight years for the browser vendors to support the tags needed to write even the most basic math. Musicians are in a similar bind, since Netscape Navigator and Internet Explorer don't support sheet music.

XML means you don't have to wait for browser vendors to catch up with what you want to do. You can invent the tags you need, when you need them, and tell the browsers how to display these tags.

## XML Describes Structure and Semantics, Not Formatting

The second thing to understand about XML is that XML markup describes a document's structure and meaning. It does not describe the formatting of the elements on the page. Formatting can be added to a document with a style sheet. The document itself only contains tags that say what is in the document, not what the document looks like.

By contrast, HTML encompasses formatting, structural, and semantic markup. `<B>` is a formatting tag that makes its content bold. `<STRONG>` is a semantic tag that means its contents are especially important. `<TD>` is a structural tag that indicates that the contents are a cell in a table. In fact, some tags can have all three kinds of meaning. An `<H1>` tag can simultaneously mean 20 point Helvetica bold, a level-1 heading, and the title of the page.

For example, in HTML a song might be described using a definition title, definition data, an unordered list, and list items. But none of these elements actually have anything to do with music. The HTML might look something like this:

```
<dt>Hot Cop
<dd> by Jacques Morali, Henri Belolo, and Victor Willis
<ul>
<li>Producer: Jacques Morali
<li>Publisher: PolyGram Records
<li>Length: 6:20
<li>Written: 1978
<li>Artist: Village People
</ul>
```

In XML the same data might be marked up like this:

```
<SONG>
  <TITLE>Hot Cop</TITLE>
  <COMPOSER>Jacques Morali</COMPOSER>
  <COMPOSER>Henri Belolo</COMPOSER>
  <COMPOSER>Victor Willis</COMPOSER>
  <PRODUCER>Jacques Morali</PRODUCER>
  <PUBLISHER>PolyGram Records</PUBLISHER>
  <LENGTH>6:20</LENGTH>
  <YEAR>1978</YEAR>
  <ARTIST>Village People</ARTIST>
</SONG>
```

Instead of generic tags like `<dt>` and `<li>`, this listing uses meaningful tags like `<SONG>`, `<TITLE>`, `<COMPOSER>`, and `<YEAR>`. This has a number of advantages, including that it's easier for a human to read the source code to determine what the author intended.

XML markup also makes it easier for non-human automated robots to locate all of the songs in the document. In HTML robots can't tell more than that an element is a `dt`. They cannot determine whether that `dt` represents a song title, a definition, or just some designer's favorite means of indenting text. In fact, a single document may well contain `dt` elements with all three meanings.

XML element names can be chosen such that they have extra meaning in additional contexts. For instance, they might be the field names of a database. XML is far more flexible and amenable to varied uses than HTML because a limited number of tags don't have to serve many different purposes.

## Why Are Developers Excited about XML?

XML makes easy many Web-development tasks that are extremely painful using only HTML, and it makes tasks that are impossible with HTML, possible. Because XML is eXtensible, developers like it for many reasons. Which ones most interest you depend on your individual needs. But once you learn XML, you're likely to discover that it's the solution to more than one problem you're already struggling with. This section investigates some of the generic uses of XML that excite developers. In Chapter 2, you'll see some of the specific applications that have already been developed with XML.

### Design of Domain-Specific Markup Languages

XML allows various professions (e.g., music, chemistry, math) to develop their own domain-specific markup languages. This allows individuals in the field to trade notes, data, and information without worrying about whether or not the person on the receiving end has the particular proprietary payware that was used to create the data. They can even send documents to people outside the profession with a reasonable confidence that the people who receive them will at least be able to view the documents.

Furthermore, the creation of markup languages for individual domains does not lead to bloatware or unnecessary complexity for those outside the profession. You may not be interested in electrical engineering diagrams, but electrical engineers are. You may not need to include sheet music in your Web pages, but composers do. XML lets the electrical engineers describe their circuits and the composers notate their scores, mostly without stepping on each other's toes. Neither field will need special support from the browser manufacturers or complicated plug-ins, as is true today.

### Self-Describing Data

Much computer data from the last 40 years is lost, not because of natural disaster or decaying backup media (though those are problems too, ones XML doesn't solve), but simply because no one bothered to document how one actually reads the data media and formats. A Lotus 1-2-3 file on a 10-year old 5.25-inch floppy disk may be irretrievable in most corporations today without a huge investment of time and resources. Data in a less-known binary format like Lotus Jazz may be gone forever.

XML is, at a basic level, an incredibly simple data format. It can be written in 100 percent pure ASCII text as well as in a few other well-defined formats. ASCII text is reasonably resistant to corruption. The removal of bytes or even large sequences of bytes does not noticeably corrupt the remaining text. This starkly contrasts with many other formats, such as compressed data or serialized Java objects where the corruption or loss of even a single byte can render the entire remainder of the file unreadable.

At a higher level, XML is self-describing. Suppose you're an information archaeologist in the 23rd century and you encounter this chunk of XML code on an old floppy disk that has survived the ravages of time:

```
<PERSON ID="p1100" SEX="M">
  <NAME>
    <GIVEN>Judson</GIVEN>
    <SURNAME> McDaniel</SURNAME>
  </NAME>
  <BIRTH>
    <DATE>21 Feb 1834</DATE> </BIRTH>
  <DEATH>
    <DATE>9 Dec 1905</DATE> </DEATH>
</PERSON>
```

Even if you're not familiar with XML, assuming you speak a reasonable facsimile of 20th century English, you've got a pretty good idea that this fragment describes a man named Judson McDaniel, who was born on February 21, 1834 and died on December 9, 1905. In fact, even with gaps in, or corruption of the data, you could probably still extract most of this information. The same could not be said for some proprietary spreadsheet or word-processor format.

Furthermore, XML is very well documented. The W3C's XML 1.0 specification and numerous paper books like this one tell you exactly how to read XML data. There are no secrets waiting to trip up the unwary.

## Interchange of Data Among Applications

Since XML is non-proprietary and easy to read and write, it's an excellent format for the interchange of data among different applications. One such format under current development is the Open Financial Exchange Format (OFX). OFX is designed to let personal finance programs like Microsoft Money and Quicken trade data. The data can be sent back and forth between programs and exchanged with banks, brokerage houses, and the like.



Cross-Reference

OFX is discussed in Chapter 2.

As noted above, XML is a non-proprietary format, not encumbered by copyright, patent, trade secret, or any other sort of intellectual property restriction. It has been designed to be extremely powerful, while at the same time being easy for both human beings and computer programs to read and write. Thus it's an obvious choice for exchange languages.

By using XML instead of a proprietary data format, you can use any tool that understands XML to work with your data. You can even use different tools for different purposes, one program to view and another to edit for instance. XML keeps you from getting locked into a particular program simply because that's what

your data is already written in, or because that program's proprietary format is all your correspondent can accept.

For example, many publishers require submissions in Microsoft Word. This means that most authors have to use Word, even if they would rather use WordPerfect or Nisus Writer. So it's extremely difficult for any other company to publish a competing word processor unless they can read and write Word files. Since doing so requires a developer to reverse-engineer the undocumented Word file format, it's a significant investment of limited time and resources. Most other word processors have a limited ability to read and write Word files, but they generally lose track of graphics, macros, styles, revision marks, and other important features. The problem is that Word's document format is undocumented, proprietary, and constantly changing. Word tends to end up winning by default, even when writers would prefer to use other, simpler programs. If a common word-processing format were developed in XML, writers could use the program of their choice.

## Structured and Integrated Data

XML is ideal for large and complex documents because the data is structured. It not only lets you specify a vocabulary that defines the elements in the document; it also lets you specify the relations between elements. For example, if you're putting together a Web page of sales contacts, you can require that every contact have a phone number and an email address. If you're inputting data for a database, you can make sure that no fields are missing. You can require that every book have an author. You can even provide default values to be used when no data is entered.

XML also provides a client-side include mechanism that integrates data from multiple sources and displays it as a single document. The data can even be rearranged on the fly. Parts of it can be shown or hidden depending on user actions. This is extremely useful when you're working with large information repositories like relational databases.

## The Life of an XML Document

XML is, at the root, a document format. It is a series of rules about what XML documents look like. There are two levels of conformity to the XML standard. The first is *well-formedness* and the second is validity. Part I of this book shows you how to write well-formed documents. Part II shows you how to write valid documents.

HTML is a document format designed for use on the Internet and inside Web browsers. XML can certainly be used for that, as this book demonstrates. However, XML is far more broadly applicable. As previously discussed, it can be used as a storage format for word processors, as a data interchange format for different programs, as a means of enforcing conformity with Intranet templates, and as a way to preserve data in a human-readable fashion.

However, like all data formats, XML needs programs and content before it's useful. So it isn't enough to only understand XML itself which is little more than a specification for what data should look like. You also need to know how XML documents are edited, how processors read XML documents and pass the information they read on to applications, and what these applications do with that data.

## Editors

XML documents are most commonly created with an editor. This may be a basic text editor like Notepad or vi that doesn't really understand XML at all. On the other hand, it may be a completely WYSIWYG editor like Adobe FrameMaker that insulates you almost completely from the details of the underlying XML format. Or it may be a structured editor like JUMBO that displays XML documents as trees. For the most part, the fancy editors aren't very useful yet, so this book concentrates on writing raw XML by hand in a text editor.

Other programs can also create XML documents. For example, later in this book, in the chapter on designing a new DTD, you'll see some XML data that came straight out of a FileMaker database. In this case, the data was first entered into the FileMaker database. Then a FileMaker calculation field converted that data to XML. In general, XML works extremely well with databases.



Cross-Reference

Specifically, you'll see this in Chapter 23, *Designing a New XML Application*.

In any case, the editor or other program creates an XML document. More often than not this document is an actual file on some computer's hard disk, but it doesn't absolutely have to be. For example, the document may be a record or a field in a database, or it may be a stream of bytes received from a network.

## Parsers and Processors

An XML parser (also known as an XML processor) reads the document and verifies that the XML it contains is well formed. It may also check that the document is valid, though this test is not required. The exact details of these tests will be covered in Part II. But assuming the document passes the tests, the processor converts the document into a tree of elements.

## Browsers and Other Tools

Finally the parser passes the tree or individual nodes of the tree to the end application. This application may be a browser like Mozilla or some other program that understands what to do with the data. If it's a browser, the data will be displayed to the user. But other programs may also receive the data. For instance, the data might be interpreted as input to a database, a series of musical notes to play, or a Java program that should be launched. XML is extremely flex-ible and can be used for many different purposes.

## The Process Summarized

To summarize, an XML document is created in an editor. The XML parser reads the document and converts it into a tree of elements. The parser passes the tree to the browser that displays it. Figure 1-1 shows this process.

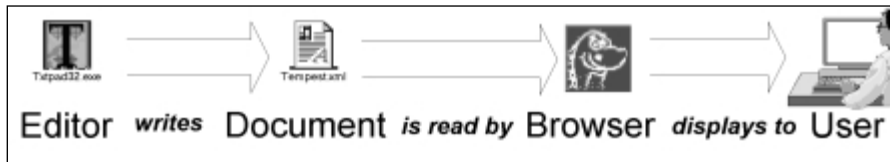


Figure 1-1: XML Document Life Cycle

It's important to note that all of these pieces are independent and decoupled from each other. The only thing that connects them all is the XML document. You can change the editor program independently of the end application. In fact you may not always know what the end application is. It may be an end user reading your work, or it may be a database sucking in data, or it may even be something that hasn't been invented yet. It may even be all of these. The document is independent of the programs that read it.



Note

HTML is also somewhat independent of the programs that read and write it, but it's really only suitable for browsing. Other uses, like database input, are outside its scope. For example, HTML does not provide a way to force an author to include certain required content, like requiring that every book have an ISBN number. In XML you *can* require this. You can even enforce the order in which particular elements appear (for example, that level-2 headers must always follow level-1 headers).

## Related Technologies

XML doesn't operate in a vacuum. Using XML as more than a data format requires interaction with a number of related technologies. These technologies include HTML for backward compatibility with legacy browsers, the CSS and XSL style-sheet languages, URLs and URIs, the XLL linking language, and the Unicode character set.

## Hypertext Markup Language

Mozilla 5.0 and Internet Explorer 5.0 are the first Web browsers to provide some (albeit incomplete) support for XML, but it takes about two years before most users have upgraded to a particular release of the software. (In 1999, my wife Beth is still

using Netscape 1.1.) So you're going to need to convert your XML content into classic HTML for some time to come.

Therefore, before you jump into XML, you should be completely comfortable with HTML. You don't need to be an absolutely snazzy graphical designer, but you should know how to link from one page to the next, how to include an image in a document, how to make text bold, and so forth. Since HTML is the most common output format of XML, the more familiar you are with HTML, the easier it will be to create the effects you want.

On the other hand, if you're accustomed to using tables or single-pixel GIFs to arrange objects on a page, or if you start to make a Web site by sketching out its appearance rather than its content, then you're going to have to unlearn some bad habits. As previously discussed, XML separates the content of a document from the appearance of the document. The content is developed first; then a format is attached to that content with a style sheet. Separating content from style is an extremely effective technique that improves both the content and the appearance of the document. Among other things, it allows authors and designers to work more independently of each other. However, it does require a different way of thinking about the design of a Web site, and perhaps even the use of different project-management techniques when multiple people are involved.

## Cascading Style Sheets

Since XML allows arbitrary tags to be included in a document, there isn't any way for the browser to know in advance how each element should be displayed. When you send a document to a user you also need to send along a style sheet that tells the browser how to format individual elements. One kind of style sheet you can use is a Cascading Style Sheet (CSS).

CSS, initially designed for HTML, defines formatting properties like font size, font family, font weight, paragraph indentation, paragraph alignment, and other styles that can be applied to particular elements. For example, CSS allows HTML documents to specify that all `H1` elements should be formatted in 32 point centered Helvetica bold. Individual styles can be applied to most HTML tags that override the browser's defaults. Multiple style sheets can be applied to a single document, and multiple styles can be applied to a single element. The styles then cascade according to a particular set of rules.

A small icon of an open book with the text "Cross-Reference" written on it.

CSS rules and properties are explored in more detail in Chapter 12, *Cascading Style Sheets Level 1*, and Chapter 13, *Cascading Style Sheets Level 2*.

It's easy to apply CSS rules to XML documents. You simply change the names of the tags you're applying the rules to. Mozilla 5.0 directly supports CSS style sheets combined with XML documents, though at present, it crashes rather too frequently.

## Extensible Style Language

The Extensible Style Language (XSL) is a more advanced style-sheet language specifically designed for use with XML documents. XSL documents are themselves well-formed XML documents.

XSL documents contain a series of rules that apply to particular patterns of XML elements. An XSL processor reads an XML document and compares what it sees to the patterns in a style sheet. When a pattern from the XSL style sheet is recognized in the XML document, the rule outputs some combination of text. Unlike cascading style sheets, this output text is somewhat arbitrary and is not limited to the input text plus formatting information.

CSS can only change the format of a particular element, and it can only do so on an element-wide basis. XSL style sheets, on the other hand, can rearrange and reorder elements. They can hide some elements and display others. Furthermore, they can choose the style to use not just based on the tag, but also on the contents and attributes of the tag, on the position of the tag in the document relative to other elements, and on a variety of other criteria.

CSS has the advantage of broader browser support. However, XSL is far more flexible and powerful, and better suited to XML documents. Furthermore, XML documents with XSL style sheets can be easily converted to HTML documents with CSS style sheets.

A small icon of an open book with the words "Cross-Reference" written on it.

XSL style sheets will be explored in great detail in Chapter 14, *XSL Transformations*, and Chapter 15, *XSL Formatting Objects*.

## URLs and URIs

XML documents can live on the Web, just like HTML and other documents. When they do, they are referred to by Uniform Resource Locators (URLs), just like HTML files. For example, at the URL <http://www.hypermedic.com/style/xml/tempest.xml> you'll find the complete text of Shakespeare's *Tempest* marked up in XML.

Although URLs are well understood and well supported, the XML specification uses the more general Uniform Resource Identifier (URI). URIs are a more general architecture for locating resources on the Internet, that focus a little more on the resource and a little less on the location. In theory, a URI can find the closest copy of a mirrored document or locate a document that has been moved from one site to another. In practice, URIs are still an area of active research, and the only kinds of URIs that are actually supported by current software are URLs.

## XLinks and XPointers

As long as XML documents are posted on the Internet, you're going to want to be able to address them and hot link between them. Standard HTML link tags can be used in XML documents, and HTML documents can link to XML documents. For example, this HTML link points to the aforementioned copy of the *Tempest* rendered in XML:

```
<a href="http://www.hypermedic.com/style/xml/tempest.xml">  
  The Tempest by Shakespeare  
</a>
```

**Note**

Whether the browser can display this document if you follow the link, depends on just how well the browser handles XML files. Most current browsers don't handle them very well.

However, XML lets you go further with XLinks for linking to documents and XPointers for addressing individual parts of a document.

XLinks enable any element to become a link, not just an A element. Furthermore, links can be bi-directional, multidirectional, or even point to multiple mirror sites from which the nearest is selected. XLinks use normal URLs to identify the site they're linking to.

**Cross-Reference**

XLinks are discussed in Chapter 16, *XLinks*.

XPointers enable links to point not just to a particular document at a particular location, but to a particular part of a particular document. An XPointer can refer to a particular element of a document, to the first, the second, or the 17th such element, to the first element that's a child of a given element, and so on. XPointers provide extremely powerful connections between documents that do not require the targeted document to contain additional markup just so its individual pieces can be linked to it.

Furthermore, unlike HTML anchors, XPointers don't just refer to a point in a document. They can point to ranges or spans. Thus an XPointer might be used to select a particular part of a document, perhaps so that it can be copied or loaded into a program.

**Cross-Reference**

XPointers are discussed in Chapter 17, *XPointers*.

## The Unicode Character Set

The Web is international, yet most of the text you'll find on it is in English. XML is starting to change that. XML provides full support for the two-byte Unicode character set, as well as its more compact representations. This character set supports almost every character commonly used in every modern script on Earth.

Unfortunately, XML alone is not enough. To read a script you need three things:

1. A character set for the script
2. A font for the character set
3. An operating system and application software that understands the character set

If you want to write in the script as well as read it, you'll also need an input method for the script. However, XML defines character references that allow you to use pure ASCII to encode characters not available in your native character set. This is sufficient for an occasional quote in Greek or Chinese, though you wouldn't want to rely on it to write a novel in another language.

A small icon of an open book with the words "Cross-Reference" written on it.

Cross-Reference

In Chapter 7, *Foreign Languages and non-Roman Text*, you'll explore how international text is represented in computers, how XML understands text, and how you can use the software you have to read and write in languages other than English.

## How the Technologies Fit Together

XML defines a grammar for tags you can use to mark up a document. An XML document is marked up with XML tags. The default encoding for XML documents is Unicode.

Among other things, an XML document may contain hypertext links to other documents and resources. These links are created according to the XLink specification. XLinks identify the documents they're linking to with URIs (in theory) or URLs (in practice). An XLink may further specify the individual part of a document it's linking to. These parts are addressed via XPointers.

If an XML document is intended to be read by human beings — and not all XML documents are — then a style sheet provides instructions about how individual elements are formatted. The style sheet may be written in any of several style-sheet languages. CSS and XSL are the two most popular style-sheet languages, though there are others including DSSSL — the Document Style Semantics and Specification Language — on which XSL is based.

**Caution**

I've outlined a lot of exciting stuff in this chapter. However, honesty compels me to tell you that I haven't discussed all of it yet. In fact, much of what I've described is the promise of XML rather than the current reality. XML has a lot of people in the software industry very excited, and a lot of programmers are working very hard to turn these dreams into reality. New software is released every day that brings us closer to XML nirvana, but this is all very new, and some of the software isn't fully cooked yet. Throughout the rest of this book, I'll be careful to point out not only what is supposed to happen, but what actually does happen. Depressingly these are all too often not the same thing. Nonetheless with a little caution you can do real work right now with XML.

## Summary

In this chapter, you have learned some of the things that XML can do for you. In particular, you have learned:

- ♦ XML is a meta-markup language that enables the creation of markup languages for particular documents and domains.
- ♦ XML tags describe the structure and semantics of a document's content, not the format of the content. The format is described in a separate style sheet.
- ♦ XML grew out of many users' frustration with the complexity of SGML and the inadequacies of HTML.
- ♦ XML documents are created in an editor, read by a parser, and displayed by a browser.
- ♦ XML on the Web rests on the foundations provided by HTML, Cascading Style Sheets, and URLs.
- ♦ Numerous supporting technologies layer on top of XML, including XSL style sheets, XLinks, and XPointers. These let you do more than you can accomplish with just CSS and URLs.
- ♦ Be careful. XML isn't completely finished. It will change and expand, and you will encounter bugs in current XML software.

In the next chapter, you'll see a number of XML applications, and learn about some ways XML is being used in the real world today. Examples include vector graphics, music notation, mathematics, chemistry, human resources, Webcasting, and more.



