



XML、 XML Schema、 XSLT 2.0和XQuery 开发详解

国内第一本
介绍XQuery
和XQJ的图书

涵盖最新的XSLT 2.0、XQuery 1.0和XQJ规范

循序渐进的讲解、
完整的代码示例、
完善的知识点演绎

孙 鑫 编著

- ◆ 国内第一本介绍XQuery和XQJ的图书
- ◆ 涵盖最新的XSLT 2.0、XQuery 1.0和XQJ规范
- ◆ 7种最常用和最新的XML技术：XML、DTD、XML名称空间、XML Schema、XPath 1.0和XPath 2.0、XSLT 1.0和XSLT 2.0、XQuery
- ◆ 4种Java解析XML文档技术：DOM、SAX、JDOM、dom4j
- ◆ 1种支持最新XQuery标准的查询API：XQJ (XQuery for Java API)



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

**Java Web
开发三部曲**

SunXin's Series
孙鑫作品系列

XML、 XML Schema、 XSLT 2.0 和 XQuery 开发详解

孙鑫 编著

电子工业出版社

Publishing House of Electronics Industry
北京·BEIJING

内 容 简 介

本书内容全面，详细讲解了目前最常用和最新的 XML 技术，包括：XML、DTD、XML 名称空间、XML Schema、XPath 1.0 和 XPath 2.0、XSLT 1.0 和 XSLT 2.0，以及 XQuery。此外，本书还介绍了如何使用 DOM、SAX、JDOM 和 dom4j 来解析和验证 XML 文档，以及使用最新的 XQJ API 来查询 XML 数据。

本书语言生动、通俗易懂、讲解细致，所有章节都提供了大量的例子，以帮助读者更好地理解所学的内容。

本书在内容的安排上独具匠心，在知识体系的讲解上由浅入深、循序渐进。

本书不仅可以作为 XML 开发的学习用书，还可以作为从事 XML 开发的程序员的参考用书和必备手册。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

图书在版编目（CIP）数据

XML、XML Schema、XSLT 2.0 和 XQuery 开发详解 / 孙鑫编著. —北京：电子工业出版社，2009.1

（Java Web 开发三部曲）

ISBN 978-7-121-07737-1

I. X… II. 孙… III. ①可扩充语言，XML—程序设计 ②可扩充语言，XSLT—程序设计
③Java 语言—程序设计 IV. TP312

中国版本图书馆 CIP 数据核字（2008）第 177417 号

责任编辑：李冰

印 刷：北京智力达印刷有限公司

装 订：北京中新伟业印刷有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：787×1092 1/16 印张：30.25 字数：610 千字

印 次：2009 年 1 月第 1 次印刷

印 数：4000 册 定价：55.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，
联系及邮购电话：(010) 88254888。

质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线：(010) 88258888。

前 言

本书内容丰富，讲解了 7 种最常用和最新的 XML 技术：XML、DTD、XML 名称空间、XML Schema、XPath 1.0 和 XPath 2.0、XSLT 1.0、XSLT 2.0 和 XQuery；4 种 Java 解析 XML 文档技术：DOM、SAX、JDOM、dom4j；1 种支持最新 XQuery 标准的查询 API：XQJ（XQuery for Java API）。

本书面向的读者

- 毫无 XML 经验的初学者
- 有一定的 XML 经验，但没有从事过 XML 开发的读者
- 正在从事 XML 开发的初中级程序员
- 正在从事 XSLT 开发的程序员
- 正在从事 XQuery 开发的程序员

本书的内容组织

全书共分 5 部分，包括了 XML 篇、XML Schema 篇、解析篇、XSLT 篇和 XQuery 篇。本书在内容的安排上独具匠心，在知识体系的讲解上由浅入深、循序渐进。

本书第 I 部分是 XML 篇，包括以下几个主题：

- XML
- DTD
- XML 名称空间

本书第 II 部分是 XML Schema 篇，详细讲解了 XML Schema 这一功能强大的模式语言。这部分内容包括以下几个主题：

- XML Schema 初窥
- 内置简单类型
- 自定义简单类型
- 复杂类型
- 派生复杂类型
- 一致性约束
- 引入其他的模式文档

本书第 III 部分是解析篇，详细讲解了 DOM、SAX、JDOM 和 dom4j 解析 XML 文档。这部分内容包括以下几个主题：

- 使用 DOM、SAX 和 JAXP 解析 XML 文档
- 使用 JDOM 解析 XML 文档
- 使用 dom4j 解析 XML 文档

➤ 解析名称空间

本书第IV部分是 XSLT 篇,介绍了 XSLT1.0/XPath 1.0, 以及最新的 XSLT 2.0/XPath 2.0。

这部分内容包括以下两个主题:

➤ XSLT 1.0

➤ XSLT 2.0

本书第 V 部分是 XQuery 篇,介绍了 XML 查询语言 XQuery, 以及在 Java 中执行 XQuery 查询的标准 API —— XQJ。这部分内容包括以下两个主题:

➤ XQuery

➤ 在 Java 中使用 XQuery —— XQJ

本书示例程序的下载

本书完整的示例程序可以从博文视点的网站 (<http://www.broadview.com.cn>) 和笔者的个人网站上下载 (<http://www.sunxin.org>)。

最后,衷心地希望本书能够给读者带来知识,带来阅读上的快乐,读者对本书的肯定就是笔者最大的欣慰。由于本书的内容较多、牵涉的技术较广,错误和疏漏之处在所难免,欢迎广大技术专家和读者指正。作者的联系方式是 csunxin@sina.com,读者也可以上作者的网站发表意见,网址是: <http://www.sunxin.org/>。

目 录

第 I 篇 XML 篇

第 1 章 XML	2
1.1 XML 的起源	2
1.2 W3C 介绍	3
1.3 关于 XML 的两个问题	4
1.4 XML 与 HTML 的比较	6
1.4.1 XML 将数据与显示分开	6
1.4.2 XML 对文档的格式要求 更加严格	7
1.4.3 XML 有且只能有一个根元素	8
1.5 XML 的编辑工具	8
1.6 XML 文档	12
1.6.1 XML 文档的结构	12
1.6.2 在 XMLSpy 中创建 XML 文档	14
1.6.3 XML 声明	15
1.6.4 文档类型声明	16
1.6.5 元素	17
1.6.6 注释	21
1.6.7 处理指令	22
1.6.8 空白处理	23
1.6.9 行尾处理	23
1.6.10 语言标识	23
1.7 格式良好的 XML	24
1.8 小结	25
第 2 章 DTD	26
2.1 在 XML 文档中引入 DTD	27
2.2 DTD 的结构	29
2.2.1 元素类型声明	30
2.2.2 实体声明	34
2.2.3 属性表声明	37
2.2.4 记号声明	46
2.3 在 XMLSpy 中创建 DTD 文档	46
2.4 有效的 XML	47
2.5 XML 处理器/解析器	48

2.6 小结	49
--------	----

第 3 章 XML 名称空间

3.1 声明名称空间	50
3.2 名称空间在元素和属性中的 运用	52
3.2.1 名称空间在元素中的运用	52
3.2.2 默认名称空间	54
3.2.3 名称空间在属性中的运用	55
3.3 名称空间和 DTD	56
3.4 小结	58

第 II 篇 XML Schema 篇

第 4 章 XML Schema 初窥	60
4.1 XML Schema 推荐标准	60
4.2 XML Schema 文档一瞥	61
4.3 XML Schema 与 DTD 的比较	62
4.4 术语明晰	64
4.5 XML Schema 概述	65
4.5.1 元素和属性的声明	65
4.5.2 元素和属性的类型	65
4.5.3 简单类型	66
4.5.4 复杂类型	68
4.5.5 全局声明和局部声明	69
4.5.6 模式与名称空间	70
4.5.7 在实例中引用模式文档	75
4.5.8 注解	78
4.6 在 XMLSpy 中创建模式文档	79
4.7 模式文档的验证	80
4.8 小结	80
第 5 章 内置简单类型	81

5.1 基于字符串的类型	83
5.1.1 string	83
5.1.2 normalizedString	83
5.1.3 token	84
5.1.4 Name	85

5.1.5 NCName	85
5.2 数字类型	86
5.2.1 float 和 double (浮点数和 双精度浮点数)	86
5.2.2 decimal (小数)	86
5.2.3 integer (整数)	87
5.3 日期和时间类型	88
5.3.1 date	88
5.3.2 time	88
5.3.3 dateTime	89
5.3.4 gYear	89
5.3.5 gYearMonth	90
5.3.6 gMonth	90
5.3.7 gMonthDay	90
5.3.8 gDay	91
5.3.9 duration	91
5.4 其他数据类型	92
5.4.1 boolean	92
5.4.2 anyURI	92
5.4.3 QName	92
5.5 小结	93
第 6 章 自定义简单类型	94
6.1 面 (facet)	95
6.1.1 限制范围	95
6.1.2 限制长度	96
6.1.3 指定精度	97
6.1.4 枚举值	98
6.1.5 模式匹配	98
6.1.6 空白处理	99
6.1.7 固定面	99
6.2 原子类型	100
6.3 列表类型	100
6.3.1 定义列表类型	100
6.3.2 限制列表类型	102
6.4 联合类型	104
6.4.1 定义联合类型	105
6.4.2 限制联合类型	106
6.5 阻止简单类型的派生	107
6.6 小结	108
第 7 章 复杂类型	109
7.1 从简单类型到复杂类型	110
7.2 xs:attribute 元素的 use、 default 和 fixed 属性	110
7.3 anyType	112
7.4 纯元素内容	112
7.4.1 sequence 组	113
7.4.2 choice 组	115
7.4.3 all 组	117
7.5 元素的出现指示符	118
7.6 元素的默认值和固定值	122
7.7 空元素	124
7.8 混合内容	125
7.9 元素组	125
7.10 属性组	127
7.11 通配符	129
7.11.1 元素通配符	129
7.11.2 属性通配符	133
7.12 小结	134
第 8 章 派生复杂类型	135
8.1 扩展派生复杂类型	136
8.1.1 扩展简单内容	136
8.1.2 扩展纯元素内容	137
8.1.3 扩展混合内容	140
8.1.4 扩展空内容	141
8.2 限制派生复杂类型	142
8.2.1 限制简单内容	142
8.2.2 限制纯元素内容	143
8.2.3 限制混合内容	145
8.2.4 限制空内容	147
8.3 在实例文档中使用派生类型	147
8.4 替换组	150
8.5 抽象元素和类型	152
8.6 控制派生类型的创建和使用	154
8.7 小结	157
第 9 章 一致性约束	158
9.1 unique 约束	159
9.2 key 约束	164
9.3 keyref 约束	165
9.4 小结	167
第 10 章 引入其他的模式文档	168
10.1 包含	168
10.2 重定义	171

10.3	导入	172
10.4	小结	174

第III篇 解析篇

第 11 章 使用 DOM、SAX 和 JAXP 解析 XML 文档 176

11.1	DOM、SAX 和 JAXP 概述	177
11.2	使用 DOM 解析 XML 文档	178
11.2.1	DOM 结构模型	178
11.2.2	DOM 树中的节点类型	180
11.2.3	DOM 解析器工厂和 DOM 解析器	186
11.2.4	JAXP 的错误类和异常类	189
11.2.5	使用 DOM 解析 XML 文档的实例	190
11.3	使用 SAX 解析 XML 文档	203
11.3.1	SAX 的处理机制	203
11.3.2	配置 SAX 解析器	207
11.3.3	SAX 解析器工厂	208
11.3.4	SAX 的异常类	210
11.3.5	ErrorHandler 接口	212
11.3.6	使用 SAX 解析 XML 文档的实例	213

11.4	使用 XML Schema 来验证 XML 文档	224
11.4.1	模式工厂和验证	224
11.4.2	与解析 API 的集成	229
11.4.3	获取类型信息	229
11.5	小结	233

第 12 章 使用 JDOM 解析 XML 文档 234

12.1	下载并配置 JDOM	235
12.2	JDOM API 介绍	235
12.3	使用 JDOM 访问 XML 文档的实例	239
12.4	小结	241

第 13 章 使用 dom4j 解析 XML 文档 242

13.1	下载并配置 dom4j	242
13.2	dom4j API 介绍	243

13.3	使用 dom4j 访问 XML 文档的实例	247
------	--------------------------------	-----

13.3.1	使用 XML Schema 来 验证 XML 文档	247
13.3.2	使用访问者模式遍历 XML 文档	249
13.3.3	使用 dom4j 的事件模型来 访问 XML 文档	251
13.4	小结	254

第 14 章 解析名称空间 255

14.1	DOM 和名称空间	256
14.2	SAX 和名称空间	259
14.3	JDOM 和名称空间	262
14.4	dom4j 和名称空间	265
14.5	小结	268

第IV篇 XSLT 篇

第 15 章 XSLT 1.0 270

15.1	XSLT 概述	271
15.2	Xalan 处理器	275
15.3	模板规则	276
15.4	<xsl:apply-templates>元素	277
15.5	<xsl:value-of>元素	278
15.6	<xsl:for-each>元素	281
15.7	匹配节点的模式	282
15.8	mode 属性	284
15.9	内置的模板规则	285
15.10	对空白的处理	287
15.11	XPath 语言	287
15.11.1	XPath 上下文	287
15.11.2	位置路径	288
15.11.3	表达式	293
15.11.4	核心函数库	295
15.12	创建结果树	300
15.12.1	创建元素和属性	300
15.12.2	创建文本	305
15.12.3	创建处理指令	307
15.12.4	创建注释	307
15.12.5	复制节点	308
15.12.6	输出格式化的数字	309

15.13	条件处理	320	16.3.1	分组	375
15.13.1	<xsl:if>	320	16.3.2	隐含文档节点（临时树）	381
15.13.2	<xsl:choose>	321	16.3.3	使用<xsl:result-document>元素输出多个文件	383
15.14	排序	322	16.3.4	<xsl:value-of>元素的改进	386
15.15	变量和参数	326	16.3.5	字符映射	386
15.15.1	变量	326	16.3.6	自定义样式表函数	389
15.15.2	参数	329	16.4	小结	390
15.16	命名模板	330			
15.17	合并样式表	331			
15.17.1	导入样式表	332			
15.17.2	包含样式表	333			
15.18	模板规则冲突的解决	333			
15.19	<xsl:output>元素	335			
15.19.1	指定输出文档的格式	336	17.1	XQuery 简介	393
15.19.2	输出 XML 声明	336	17.2	查看 XQuery 的查询结果	394
15.19.3	输出文档类型定义	337	17.2.1	XMLSpy 和 Stylus Studio	394
15.19.4	输出 CDATA 段	338	17.2.2	Saxon	396
15.19.5	指定文档缩进	340	17.2.3	DataDirect XQuery	397
15.19.6	指定媒体类型	340	17.3	XQuery 基础	398
15.20	XSLT 中的函数	340	17.3.1	处理模型	398
15.21	数字格式化	341	17.3.2	表达式上下文	398
15.22	查询和分组	343	17.3.3	数据模型	399
15.23	处理多个输入文档	352	17.3.4	类型	400
15.24	JAXP 中的 XSLT API	356	17.3.5	注释	400
15.24.1	转换器工厂	356	17.4	表达式	401
15.24.2	Transformer 和 Templates	356	17.4.1	基本表达式	401
15.24.3	一个实例	358	17.4.2	路径表达式	402
15.25	在 XMLSpy 中创建样式表		17.4.3	序列表达式	402
	文档	360	17.4.4	算术表达式	403
15.26	小结	362	17.4.5	比较表达式	404
第 16 章	XSLT 2.0	363	17.4.6	逻辑表达式	405
16.1	Saxon 处理器	364	17.4.7	FLWOR 表达式	406
16.2	XPath 2.0	365	17.4.8	有序和无序表达式	410
16.2.1	一切都是序列	365	17.4.9	条件表达式	411
16.2.2	for 表达式	366	17.4.10	量化表达式	412
16.2.3	条件表达式	367	17.4.11	作用于序列类型的表达式	413
16.2.4	限定性表达式	369	17.4.12	验证表达式	417
16.2.5	类型	371	17.4.13	扩展表达式	418
16.2.6	日期和时间	372	17.5	查询的结构	419
16.2.7	函数	373	17.5.1	主模块和库模块	419
16.3	XSLT 2.0 的新特性	375	17.5.2	版本声明	421

17.5.3	序言	421
17.5.4	设置器	422
17.5.5	名称空间声明	423
17.5.6	默认名称空间声明	424
17.5.7	模式导入	425
17.5.8	模块导入	426
17.5.9	变量声明	427
17.5.10	函数声明	430
17.5.11	选项声明	431
17.6	小结	432

第 18 章 在 Java 中使用 XQuery——XQJ		433
18.1	XQJ 简介	433
18.2	开发一个 XQJ 应用	434
18.3	对 XQuery 上下文的支持	439
18.4	映射 XQuery 数据模型	440
18.5	对 XQuery 类型系统的支持	441
18.6	XQMetaData 接口	441
18.7	小结	442
附录 A 快速掌握 HTML		443



第 I 篇

XML 篇

XML 已经成为大家耳熟能详的一个词汇，话说微软力推的少数几项非自己专有技术的标准中，就有大名鼎鼎的 XML。XML 一经推出，就得到了业界 IT 巨头的响应，一时间应者云集，XML 很快就在各行各业中展露了它的身影。XML 能够独立于计算机平台、操作系统和编程语言来表示数据，凭借其简单性、可扩展性、交互性和灵活性在计算机工业中获得了世界范围的支持和采纳。人们使用 XML 作为配置文件的首选格式，使用 XML 在不同的系统、不同的数据库，以及不同的软件之间传输数据，使用 XML 作为简易数据库来保存企业信息数据。

XML 的流行使得很多软件开发人员在开发项目的时候，不管是否需要，都在项目中使用了 XML，从而导致软件系统性能的降低，这种盲目追求新技术或流行技术的做法并不可取。

为了帮助读者对 XML 技术有一个清晰而全面的认识，本篇将重点讲述 XML、DTD，以及 XML 名称空间。本篇的内容包括：

- ◆ 第 1 章 XML
- ◆ 第 2 章 DTD
- ◆ 第 3 章 XML 名称空间



第 1 章

XML

本章要点

- 了解 XML 的起源
- 弄清楚 XML 与 HTML 的区别
- 学会使用 XMLSpy
- 了解 XML 的物理结构
- 掌握 XML 的逻辑结构
- 学会编写格式良好的 XML 文档

XML 对我们来说已不再陌生，其相关概念和知识在网络上随处可见，有关 XML 的应用也越来越多。本章的目的是帮助读者快速了解和掌握 XML，为后面章节的学习打下基础。

1.1 XML 的起源

XML 的全称是 Extensible Markup Language，意思是可扩展的标记语言，它是标准通用标记语言（Standard Generalized Markup Language，SGML）的一个子集。那 SGML 又是什么呢？

在 20 世纪 80 年代早期，IBM 提出在各文档之间共享一些相似的属性，例如字体大小和版面。IBM 设计了一种文档系统，通过在文档中添加标记，来标识文档中的各种元素，IBM 把这种标识语言称做通用标记语言（Generalized Markup Language，GML）。经过多年的发展，1984 年国际标准化组织（ISO）开始对此提案进行讨论，并于 1986 年正式发布了为生成标准化文档而定义的标记语言标准（ISO 8879），称为新的语言 SGML，即标准通用标记语言。

SGML 功能非常强大，是可以定义标记语言的元语言，然而由于 SGML 过于复杂，不适合在 Web 上应用，因此 W3C 组织在 1996 年便开始设计一种可扩展的标记语言，以便

能将 SGML 的丰富功能与 HTML 的易用性结合到 Web 应用中。1998 年 2 月 10 日，W3C 发布了 XML 1.0 标准，其目的是为了在 Web 上能以现有的超文本标记语言（HTML）的使用方式提供、接收和处理通用的 SGML。XML 是 SGML 的一个简化子集，它以一种开放的、自我描述的方式定义了数据结构。在描述数据内容的同时能突出对结构的描述，从而体现出数据与数据之间的关系。

本书介绍的 XML 主要遵循 W3C 于 2006 年 8 月 16 日发布的 XML 1.0 推荐标准的第四版，读者可以在 <http://www.w3.org/TR/REC-xml/> 了解到此标准的详细内容。

W3C 组织于 2004 年 2 月 4 日，发布了 XML 1.1 的推荐标准，并于 2006 年 8 月 16 日发布了 XML 1.1 推荐标准的第二版，这是最新的 XML 版本，不过由于目前大多数的应用还是基于 XML 1.0 的推荐标准，因此本书也将遵照 XML 1.0 规范来讲述。如果读者想要了解 XML1.1 规范的内容，可以参看网址：<http://www.w3.org/TR/2006/REC-xml11-20060816/>。

1.2 W3C 介绍

W3C 是万维网联盟（World Wide Web Consortium）英文的缩写，成立于 1994 年 10 月，它以开放论坛的方式来促进开发互通技术（包括规格、指南、软件和工具），以激发网络的全部潜能。万维网联盟（W3C）的会员包括软件开发商、内容提供商、企业用户、通信公司、研究机构、研究实验室、标准化团体，以及政府，会员中的一些知名 IT 企业包括：IBM、Microsoft、America Online、Apple、Adobe、Macromedia、Sun Microsystems 等。

W3C 自成立以来，已发布了 100 多份 Web 技术规范，领导着 Web 技术向前发展。

W3C 认为自身不是官方组织，因此将它正式发布的规范称为**推荐（建议）标准**，意思是进一步标准化的建议，但是由于该组织自身的权威性往往成为事实上的标准。一项技术要成为 W3C 的推荐标准，需要经过以下 7 个步骤。

（1）W3C 收到提交（W3C Submissions）

任何 W3C 的成员都可以向联盟提交希望成为 Web 标准的某项建议。如果建议的内容在 W3C 的工作范围内，W3C 将决定是否要对此开展工作。

（2）W3C 发布注释（W3C Notes）

通常，一项对 W3C 的提交会成为一份注释。注释作为一份公共的文档，是对建议的描述。

注释仅供讨论使用。一项注释的发布，并不表示 W3C 对其认可了。注释的内容由提交此注释的会员来编辑，而不是 W3C。注释可以在发布后随时被更新、替换或废弃。注释的发布，并不表示 W3C 已经开始任何与注释相关的工作。

（3）W3C 成立工作组（W3C Working Groups）

当一项提交被 W3C 认可之后，就成立由成员和一些对此感兴趣的团体参加的工作组。通常情况下，工作组将定义一个时间表，并发布一个提议标准的工作草案，描述当前的工作进展。



(4) W3C 发布工作草案 (W3C Working Drafts)

W3C 通常会在其网站上 (<http://www.w3.org>) 发布工作草案，以及一个公众讨论的邀请。工作草案会说明当前的工作进展。由于工作草案的内容可随时被更新、替换或废弃，所以不应把它作为工作的依据。

(5) W3C 发布候选推荐标准 (W3C Candidate Recommendations)

当规范比较复杂时，可能需要成员和软件提供商花费更多的时间来试用或测试。有时候，这些规范以候选推荐标准的形式发布。它与工作草案一样都是进展中的工作文件，所以不应把它作为工作的依据。该文档随时可以被更新、替换或废弃。

(6) W3C 发布提议的推荐标准

提议的推荐标准的发布标志着工作组的工作到了最后阶段。提议的推荐标准仍然是一个进展中的工作文件，仍然可以被更新、替换或废弃。虽然提议的推荐标准还没有被 W3C 正式认可，但是它在内容和时间上离最终的推荐标准已经非常接近了。

(7) W3C 发布推荐标准

推荐标准经过 W3C 的成员审阅，并由 W3C 的主任加盖正式批准图章，而最终成为规范。W3C 的推荐标准是一个稳定的文档，可以作为工作中的参考资料。

1.3 关于 XML 的两个问题

1. XML 是 HTML 的扩展吗

HTML 的全称是 Hypertext Markup Language (超文本标记语言)，而 XML 的全称是 eXtensible Markup Language (可扩展的标记语言)，这很容易让人联想到 XML 是通过增加新标记来扩展 HTML 的一种标记语言。实际上 HTML 和 XML 在标记语言中处于不同的层次。

下面我们通过 HTML 文档和 XML 文档的对比（分别如例 1-1 和例 1-2 所示），来弄清楚 XML 的一些概念。

例 1-1 HelloWorld.html

```
<html>
  <head>
    <title>这是一个欢迎的例子</title>
  </head>
  <body>
    你好！欢迎你！
  </body>
</html>
```

例 1-2 HelloWorld.xml

```
<?xml version="1.0" encoding="GB2312"?>
<欢迎词>
    <标题>这是一个欢迎的例子</标题>
    <内容>你好！欢迎你！</内容>
</欢迎词>
```

从上面的两个例子，我们可以看出：

① 在编写 HTML 文档时，所有的标记都已经固定下来（如<html>、<body>等），我们不能去创造新的标记；而在编写 XML 文档时，我们可以任意地创建新的标记，包括中文的标记（如<欢迎词>、<内容>等）。所以说 XML 是可扩展的标记语言。

② 在编写 XML 文档时，没有一套标准的标记供我们选择使用，需要我们自己去创建标记，所以我们说 XML 是创建标记语言的元语言。



提 示 XML 在设计之初，就考虑到了国际化的问题，同 HTML 4.01 一样，XML 也是基于 ISO/IEC 10646 字符集标准中定义的通用字符集（Universal Character Set, UCS），等价于 Unicode 2.0。

2. SGML、HTML 和 XML 之间是什么关系

SGML 是一种在 Web 发明之前就早已存在的使用标记来描述文档资料的通用语言。它是一种定义标记语言的元语言。HTML 和 XML 都是从 SGML 发展而来的标记语言，因此，它们有一些共同点，如相似的语法和标记的使用。不过 HTML 是在 SGML 定义下的一个描述性的语言，只是 SGML 的一个应用，其 DTD（参见第 2 章）作为标准被固定下来，而 XML 是 SGML 的一个简化版本，是 SGML 的一个子集，严格意义上来说，XML 仍然是 SGML。

HTML 不能用来定义新的应用，而 XML 可以，例如 RDF 和 CDF 都是使用 XML 定义的应用。事实上，XML 和 SGML 是兼容的，但又没有 SGML 那么复杂，它被设计用于有限带宽的网络，如 Internet。XML 规范的制定者之一 Tim Bray 说，“XML 的设计出发点是取 SGML 的优点，去除复杂的部分，使其保持轻巧，可以在 Web 上工作”。

下面我们用一个比喻来描述它们三者之间的关系。假如我们制定了一套设计图形元素的规则，称为 SGML。现在为了规范建筑行业建筑图的绘制，需要使用一套固定的图形元素，于是我们根据 SGML 这种设计图形元素的规则设计了一套专门用于绘制建筑图的图形元素，那么这套图形元素就称为 HTML。然而由于 SGML 过于复杂，不利于在其他行业推广，因此我们对 SGML 进行了简化，重新制定了一套设计图形元素的规则，这就是 XML。由于 XML 是 SGML 的一个子集，因此它也可以用来设计图形元素。

HTML、SGML 和 XML 将继续用于其适合的地方，它们中的任何一个都不会使其他一个废弃。对于像新闻、网络日记、论坛留言等大部分短期的数据，HTML 仍是在 Web 上快速发布数据的最简单的方法。如果数据要长期使用，并且需要更多的一些结构，我们更推荐使用 XML。不同于 HTML 和 XML，SGML 可能永远不会在 Internet 上被广泛接受，因为它不是为某个网络协议而设计的，也从来没有为某个网络协议的需求而优化过。对于高端的、复杂结构的发布应用，SGML 将继续使用。



1.4 XML 与 HTML 的比较

通过 XML 与 HTML 的比较，我们能够更好地理解和掌握 XML 的优点。

1.4.1 XML 将数据与显示分开

我们看附录中的一个例子。

例 1-3 SecondPage5.htm

```
<html>
  <head>
    <title>静夜思</title>
  </head>
  <body>
    <center>
      <h2><font color="red">静夜思</font></h2>
      <b>作者：李白</b>
      <hr color="blue">
      <p>
        <b><i><font size=3 color="green">床前明月光，疑是地上霜。<br>
          举头望明月，低头思故乡。</font></i></b>
      </p>
    </center>
  </body>
</html>
```

在上面的这个 HTML 文档中，“静夜思；作者：李白；床前明月光，疑是地上霜。举头望明月，低头思故乡。”这些是要显示的数据，`<center>`表示让内容居中显示，`<h2>`表示 2 号标题，``表示以粗体显示，`<i>`表示以斜体显示。

可以看到，HTML 文档将数据，页面的排版，以及页面的表现形式混合在了一起，如果要增加新的数据，相应地就要调整数据的排版与显示方式。当我们从其他的地方（例如，数据库、文本文件）得到数据后，一旦放入 HTML 文档中，整个数据就会被打乱，数据和 HTML 标记混合在了一起，数据本身就得无法辨析了。另一方面，随着电子商务等网络应用的流行，不同系统、不同平台、不同软件的信息交换日益频繁，而 HTML 本身的这些限制，导致了 HTML 在日益广泛的 Web 应用中，显得捉襟见肘，为此，就有了 XML 的产生。

我们将例 1-3 中的 HTML 文档中的数据部分改为用 XML 文档来表示，如例 1-4 所示。

例 1-4 SecondPage5.xml

```
<?xml version="1.0" encoding="GB2312"?>
<poem>
  <title>静夜思</title>
  <author>李白</author>
  <content>
```

```
<line>床前明月光</line>
<line>疑是地上霜</line>
<line>举头望明月</line>
<line>低头思故乡</line>
</content>
</poem>
```

在这个文档中，我们用<title>（此处的<title>是我们自己创建的标记）指明标题是“静夜思”，用<author>表明作者是“李白”，而内容是“床前明月光，疑是地上霜。举头望明月，低头思故乡”这四句诗。那么这个 XML 文档将如何显示呢？从文档中我们看不出这些数据将如何显示，这是因为：XML 文档不能描述页面的排版和表现形式，它只是用于描述数据和数据的结构。也就是说，XML 将数据和显示分开了，我们可以为这些数据设计不同的排版和表现形式，而数据本身不需要做任何的修改。

采用 XML 来表示数据，我们能够很容易地读懂 XML 文档，而计算机也能够很好地进行识别和处理。XML 表示数据的方式真正做到了独立于应用系统，并且数据能够重用，一份数据可以应用于不同的场合。有时候，XML 文档也被看作是文档的数据库化和数据的文档化。

1.4.2 XML 对文档的格式要求更加严格

由于 HTML 文档格式非常松散，导致了 HTML 文档解析的复杂性，也造成了浏览器兼容的问题，所以 XML 从一开始，就对文档的格式制定了非常严格的标准，凡是符合这一标准的 XML 文档就是格式良好的 XML 文档（Well-Formed XML Documents）。

（1）开始标签必须要有一个配套的结束标签

在 HTML 文档中，可以直接使用<p>、<tr>、<td>等标签，而不用加结束标签，在 XML 中，开始标签和结束标签必须配套，也就是必须写成<p>…</p>、<tr>…</tr>或<td>…</td>。

（2）空元素标签必须被关闭

在 HTML 文档中，可以使用
、<hr>、等单标签，而在 XML 中，空元素标签必须被关闭。空元素标签采用斜杠（/）来关闭，例如：
、<hr/>、。

（3）所有的标签都区分大小写

在 HTML 文档中，标签是不区分大小写的，<tr>和</TR>是 tr 元素的开始标签和结束标签，但是在 XML 中，<tr>和<TR>是两个不同的标签，开始标签和结束标签的大小写形式必须一致。

（4）所有的标签都必须合理嵌套

在 HTML 文档中，<i>…</i>是允许的，但是在 XML 中，这是错误的。在 XML 中，所有的标签都要成对出现，合理嵌套，正确的形式是：<i>…</i>。