

计算机技术译林精选系列

XML & SGML 参考手册

Rick Jelliffe 著

潇湘工作室 译

本书附盘可从本馆主页 <http://lib.szu.edu.cn/>
上由“馆藏检索”该书详细信息后下载，
也可到视听部复制

人民邮电出版社

计算机技术译林精选系列
XML & SGML 参考手册

- ◆ 著 Rick Jelliffe
 译 潇湘工作室
 责任编辑 李 际
- ◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街 14 号
 邮编 100061 电子函件 315@ pptph.com.cn
 网址 <http://www.pptph.com.cn>
 北京汉魂图文设计有限公司制作
 北京顺义振华印刷厂印刷
 新华书店总店北京发行所经销
- ◆ 开本:787 × 1092 1/16
 印张:33
 字数:818 千字 2000 年 10 月第 1 版
 印数:1 - 5 000 册 2000 年 10 月北京第 1 次印刷

著作权合同登记 图字:01 - 1999 - 0430 号

ISBN 7-115-08727-X/TP·1778

定价:65.00 元

版权声明

Rick Jelliffe: The XML & SGML Cookbook

Authorized translation from the English language edition
published by Prentice Hall PTR.

Copyright © 1998 by Rick Jelliffe.

All rights reserved. No part of the book may be reproduced
or transmitted in any form or by any means, electronic or
mechanical, including photocopying, recording or by any
information storage retrieval system, without permission
from the Publishers.

Chinese Simplified language edition published by People's
Posts & Telecommunications Publishing House.

本书英文版由 Prentice Hall PTR 出版。人民邮电出版社取得授权翻译出版中文简体版。

未经出版者许可，对本书任何部分不得以任何方式
或任何手段复制和传播。

版权所有，侵权必究。

译者序

1998年2月，W3C正式批准了XML（扩展标记语言）的标准定义，XML可以在文本文档中标记结构，也就是说，它可以对文档和数据进行结构化处理，从而使它们能在部门、客户和供应商之间进行交换。现在，尽管XML尚未广泛使用，但由于它能让公司共享和使用分布在不同平台上的信息，因而它正在引起人们越来越多的注意。

利用XML，可以实现动态内容生成、企业集成和应用开发。现在，人们还要支持过时的旧系统，所以需要合并多种旧式系统中的数据，在这方面，XML可以为我们减轻负担。

本书是介绍XML系列书中的一本，本套书共有3本，分别是《XML用户手册》、《XML & SGML参考手册》和《XML应用实例——建立电子商务应用》。本套书适用于利用XML进行各种应用（如网上银行、“推”技术、Web自动化操作、数据库发布、软件销售等）的人员。本套书的主要内容有：

- 在《XML用户手册》中，介绍XML的定义和作用，它与HTML的区别以及各种概念。在这里可以学到在现实世界中成功的XML应用以及富有经验的开发人员取得成功的秘诀，了解他们是如何应用XML来进行电子商务、内容管理、结构化、创建和表达的。
- 在《XML应用实例——建立电子商务应用》中，循序渐进地介绍如何建立完整的XML电子商务应用程序。读者阅读这本书，可以理解DOM的关键作用，并可以看到XML在高级电子商务应用程序中的成功案例。
- 在《XML & SGML参考手册》中，给出了在各种常见的编辑结构中应用的元素，它可以作为参考资料使用。另外，还给出了XML、SGML、HTML、TEI、CALS和国际化应用程序的实用技巧、警告信息以及解决方案。

在本书中，有建立XML Web应用的各种技巧、代码、工具和资源，它是一本很全面的参考资料。

参加本书翻译的人员主要有胡斌辉、彭少熙、王晓鹏、陶涛、滕杰、肖展业、孟丽艳等，另外，孙俊峰、贺民、龚亚平等也做了大量的工作。全书最后由贺军统稿。

由于水平所限，翻译错误和疏漏之处在所难免，请不吝指正。

潇湘工作室 2000年6月

原书序

每位厨师都有其独到的烹饪方法。不必说，那就是可以增加菜肴滋味的酱油、高汤和调料的某种组合。

Rick Jelliffe 也不例外。他是非常棒的“厨艺”大师，他烹饪的是结构化的信息，也就是支持企业运转，以及支撑着电子商务和 World Wide Web 的数据和文档。10 年来，他一直在帮助世界各地的客户掌握结构并赢得对自己的信息处理的控制。

Rick 的秘密配料就是 SGML（用于结构化信息的国际性标准）和 XML（为 World Wide Web 设计的 SGML 的子集）。事实上，Rick 一直在开发这些标准的委员会中与我共事。

多年以来，Rick 收集了许多令人惊异的“烹饪”秘诀，并将在本书中与你分享。当然，其中包括了所有的基础知识，但在诸如东亚语言和国际化这样的主题中也有相当深奥的专家经验。

Rick 不仅是个“美食家”，而且还是个善于讲故事的人。为了帮助理解和保持概念，书中还点缀着迷人的历史和文化背景。

所以，请读一读本书，和 Rick 一起加入到结构化信息技术的盛宴中来。或者，如同他说的那样：

“Come on in and we'll throw another structure on the barbie!”
(快点，我们来烤另一种结构)。

Charles F. Goldfarb

Saratoga

加利福尼亚

前　　言

当我今天完成这本书时，我从网上获得了两份激动人心的新出版物，新鲜得甚至还留着激光打印机的余温：一份是 XML 1.0 的最终文本（World Wide Web 联盟的 SGML 子集，即 Extensible Markup Language（扩展标记语言）），另一份是 Web SGML（对 SGML 标准——ISO 8879——的修正和改进，像 XML 一样用于 WWW）。

这两份出版物共同掀起了电子出版业的革命，它们使得从前大型团体出版才可获得的 SGML 的优势现在在桌面上也可以获得了。最重要的是，它们令你——数据的拥有者和数据系统的创建者——可以控制自己的文档：SGML 及其针对 Web 优化的子集 XML 代表了 Open Systems（开放系统）运动的胜利。我们不仅能够获得用于网络协议和语言的开放系统，还能够获得激活它们的数据所用的开放系统。

本书是这个全新世界的一本实践指南——SGML 和 XML 元素类型集合和实体集合的构想和声明，这些集合能实现重要而有用的文档结构。我尝试着放入许多资源，这些资源用其他方法很难搜寻到。并且我特别注意了设法说清楚 XML 中更深入的模型，该模型对从 HTML 或其他专有标记语言转来的读者来说可能会觉得陌生：处理指令和特殊的符号。

这不是一本语法指南：指南类的书可以找到很多不错的，包括本系列中先前的几本。我特别留心避免引用 SGML 的声明，避免引用 SGML 中那些很少实现的、可选的功能。我把对一个领域的特别详细的论述包括了进来，该领域至今为止还几乎没有被详细地论述过：字符、字形和国际化。

一、规律、结构、模式和格式

本书是针对规律、结构、模式和格式的。在本书中：

- 规律（order）是指潜在的、抽象的（有时是难以言说的）关系和事物的本质。
- 结构（structure）是指如何用有形的标记来捕获规律。
- 模式（pattern）是指用于创建结构的一类模式或方法（例如，制衣意义上的“模板”不同于文字处理意义上的“模式”）。
- 格式（form）指一个结构和另一个结构之间特殊的一致性（例如，具

体事物层面意义上的“格式”不同于形而上学意义上的“格式”）。

和所有的好食谱一样，除了说明“如何”生成结构外，本书还会设法解释“为什么”，会探究有没有别的选择，会给出各种正面和反面的意见。

像 SGML 这样非常通用的技术的自由度可能会令新文档系统的设计者感到不适。能够向任何方向移动并不是那么舒服的，如果你不能容忍走错方向的可能！幸运的是，在 SGML 作为国际性标准存在的 10 年中，适用于许多常用的文档模板的集成方法和解决方案已经显现出来。本书尝试列出并讨论它们中最好的、最有教益的和最有用的方法。

二、文档系统

对系统的局限性和除了文档类型声明之外的各种因素的考虑常常成为决定工程是否成功的关键。正因为如此，本书更关注于“文档系统设计者”（document-system designer）而不仅仅是“DTD 编写者”（DTD writer）。我们不可能在真空中管理信息，文档是作为文档系统的一部分而存在的。有时系统是封闭的，有时则是开放的。如果这些信息有足够的价值，可以为管理工作提供保证，那么在创建一个大型的 DTD 时，为整个文档系统做考虑就会很有用。

无论如何，本书对制作 XML 文档的人们以及那些甚至从未把 SGML 当作符号来创建正式的元素类型声明的人们也有用处。如果你是这些人中的一员，我强烈推荐你学习并至少尝试着在非正式的文档中使用 SGML 内容模型：SGML 可提供非常方便并已得到好评的符号，它们适用于许多种类的结构，你还可获得图形化和可视化工具的帮助。

所以，即使可以在文档的任何地方找到规律，许多情况下结构也是松散的，其中包含着例外或者是不完全的结构。因此，本书中用于元素类型集合的模板可以视为原型和标本，你可以利用并改造它们来满足自己的特殊需要，而不是那种你只能顺从地剪切和粘贴的模板。

文档系统设计者需要知道 DTD 简洁性的局限。设计者认为可以作为著者原型的模板可能实际上并没有太大的用处。只有在能够揭示出一些实际规律，而决不是强加给一种似是而非的规律时，模板的使用才是成功的。

文档系统设计者往往具有拒绝混乱、爱好整洁的癖好，有时会为想在没有规律的地方找到规律而付出代价：那只是从先前的文档中得到的“幻象”。所以，本书除了给出各种模板之外，也为选择模板提供了一些原则。必须用成功的欲望来节制对简洁的要求。我希望那些翻开这本书并期望能够得到快刀斩乱麻般干净利落的解决方案的读者能提高自己的能力，你会理解最适合于自己的个人需求的那些问题并权衡利弊。

三、术语

XML 正在把大批人员从各种学科和技术领域带入 SGML 的世界，所以，目前术语还有相当大的变异和重复。为了控制叙述的方式，我使用了一些常见的经过简化的术语，它们着重于已使用的 SGML 关键字，如下表所示：

本书	ISO-ese
ANY 元素	已声明的内容类型为 ANY 的元素
CDATA 属性	已声明的值为 CDATA 的属性
CDATA 元素	已声明的内容类型为 CDATA 的元素
CDATA 实体	CDATA 实体
CDATA 标记段	CDATA 标记段
容器元素	具有子元素的元素
EMPTY 元素	空元素
ID 属性	已声明的值为 ID 的属性
IDREF 属性	已声明的值为 IDREF 的属性
NCDATA 实体	NCDATA 实体
NMTOKEN 属性	已声明的值为 NMTOKEN 的属性
RCDATA 元素	已声明的内容类型为 RCDATA 的元素
SDATA 实体	SDATA 实体
SUBDOC 实体	子文档实体

在本书中，“一个属性 ID”是指“一个名称为 ID 的属性”；“一个 ID 属性”是指“一个已声明的值为 ID 的属性”；但“该属性 ID”是指“示例代码段中名为 ID 的属性”。好的用法是“一个属性 ID”应该是“一个 ID 属性”，同样，“一个属性 IDREF”最好也是“一个 IDREF 属性”。

我打算在 Prentice Hall PTR 的 Web 站点 www.phptr.com 上维护一个 Web 网页，内容是本书的勘误表。

Rick Jelliffe
澳大利亚，悉尼

目 录

第一部分 文档系统

第 1 章 文档和出版物	3
1.1 明确的和隐含的文档类型.....	4
1.2 热爱、臃肿和明智.....	4
1.3 出版物的六视图模型.....	5
1.3.1 查看版面设计	6
1.3.2 查看页面对象	7
1.3.3 查看字形	7
1.3.4 查看字符	8
1.3.5 查看编辑性结构	9
1.3.6 查看主题结构	10
1.3.7 相关性流程	11
1.4 时尚、趋势、争论.....	14
1.5 HTML 和 SMDL	16
1.6 关于非文本.....	16
1.7 文档与 API 的比较	18
1.8 小结.....	19
第 2 章 标记的本质	20
2.1 什么是好的标记.....	21
2.2 普通标记和特殊标记.....	23
2.3 哪一个更好：通用标记还是特殊标记.....	25
2.4 基本形式.....	26
2.5 嵌入其他类型的数据.....	27
2.6 世界上最糟糕的 DTD	29
2.7 小结.....	30

第3章 软件工程	31
3.1 DTD与模板	32
3.1.1 可重用组件	32
3.1.2 结构格式	32
3.1.3 信息单元	33
3.1.4 结合与耦合	34
3.2 瀑布与螺旋方式	35
3.2.1 图表	36
3.2.2 Maler 和 el Andoloussi 的方法论	38
3.3 考察与原型化	39
3.3.1 原型化	40
3.3.2 考察性的 DTD 设计	40
3.4 人类的角度	41
3.4.1 视角分析	41
3.4.2 方案分析	43
3.4.3 用户界面是文档	44
3.4.4 涉及因素	45
3.4.5 有用的技巧	45
3.5 小结	46
第4章 实现的选择	47
4.1 DTD样式校验表	47
4.2 是否需要完整的 SGML	54
4.3 准 SGML	55
4.3.1 非标准的通用标记	55
4.3.2 HTML	56
4.3.3 XML	57
4.3.4 属于自己的简化 SGML	58
4.3.5 带有用户扩展的 SGML	61
4.4 经验法则	62
4.4.1 语言类比	62
4.4.2 对象关系	63
4.4.3 事件	63
4.4.4 顺序	64
4.5 小结	64
第5章 使用中的文档	66
5.1 仅有声明还不够	66

5.2 处理 SGML	68
5.2.1 SGML 工具	68
5.2.2 文本工具	69
5.2.3 存储管理工具	71
5.2.4 混合工具	72
5.3 SGML 的常规步骤	73
5.3.1 “擦洗”	73
5.3.2 “按摩”	75
5.3.3 “拧干”	75
5.4 DTD 的成长	76
5.5 DTD 失败的 10 大原因	78
5.6 小结	81
5.7 设计原则	81

第二部分 文档模式

第 6 章 常用属性	85
6.1 SGML	85
6.2 HTML&XML	86
6.3 XLL	86
6.4 TEI	87
6.5 SGML 扩展机制	88
6.5.1 默认值列表	88
6.5.2 用于元素的数据属性	89
6.5.3 限制 IDREF 的目标元素类型	90
6.5.4 常用的数据属性	90
6.6 未指定的属性	91
6.7 去掉修饰	91
6.8 结构化形式	92
第 7 章 文档外壳	94
7.1 HTML	94
7.2 信息单元	95
7.3 简单开头的优点	97
第 8 章 段落	98
8.1 段落与文本块的比较	98
8.2 段落与段落组的比较	100

8.3 段落内容	100
8.4 段落内嵌套的段落	101
8.5 小段	102
8.6 ID 属性	104
8.7 段落的发展	104
8.8 段落分界符	105
8.9 再谈段落组	106
第 9 章 序列	108
9.1 序列示例	108
9.2 糟糕的混合内容	110
9.3 简化线性形式	111
第 10 章 命名数据	112
10.1 带字段的文本	112
10.1.1 带字段文本的序列	113
10.2 元素引用	115
10.3 描述表	117
10.4 导入 ASCII 转储	118
10.5 使用参数实体进行模式和类型扩展	119
第 11 章 表格	121
11.1 直接标记与元素引用的比较	121
11.2 简单的 HTML 样式表	122
11.3 ICADD 表格	122
11.4 CALS 表格	123
11.5 HTML 4 表格	125
第 12 章 交互式系统	126
12.1 实体	127
12.2 元素	127
12.3 处理指令	129
第 13 章 正式公共标识符	131
13.1 SGML 开放实体编目	132
13.2 SGML 和 MIME	133

第 14 章 数据内容表示法	135
14.1 表示法	135
14.2 表示法的一些 FPI	137
14.2.1 ISO 标准	137
14.2.2 时间和空间	146
14.2.3 非标准	147
第 15 章 正式系统标识符	151
15.1 正式系统标识符	151
15.2 成为 FSI 用户	153
第 16 章 嵌入式表示法	156
16.1 命名	156
16.2 样式表和脚本	157
16.3 定义数据类型	159
16.3.1 使用标准表示法名称的词法分类	159
16.3.2 词法模型的词法分类	160
16.3.3 使用 HyLex 的日期属性	161
16.3.4 日期使用 POSIX 常规表达式	162
16.4 嵌入其他表示法	162
16.5 分段交换	163
第 17 章 组织和记录 DTD	164
17.1 核心元素类型集	164
17.2 基础和派生 DTD	165
17.3 结构形式	167
17.4 DTD 版本	167
17.5 多遍 DTD	168
17.6 未说明的元素	168
17.6.1 简单的形式	168
17.6.2 更丰富的形式	169
17.7 记录 DTD	170
17.7.1 外部文档	170
17.7.2 注释	170
17.7.3 额外的要求	171
17.7.4 文档实例中的描述	172

第三部分 字符和字形

第 18 章 字符和字形	175
18.1 ISO 字符/字形模型	175
18.2 演变历史	177
18.3 现代印刷脚本	179
18.4 字符指令表	181
18.4.1 使用实体	181
18.4.2 使用元素	182
18.4.3 字符	183
18.5 整理	183
18.5.1 英语的简单整理	183
18.5.2 西欧语言的整理	183
18.5.3 模糊转换	184
18.5.4 显式标记	184
18.6 小结	185
第 19 章 字体、脚本和语言	186
19.1 字体	186
19.1.1 西方	186
19.1.2 东方	187
19.1.3 指定确切的字体	188
19.1.4 设计组	189
19.2 脚本代码	195
19.3 语言代码	202
19.4 国家/地区代码	216
19.5 多语种文档	223
19.5.1 内联地方化	223
19.5.2 实体	224
19.5.3 元素和属性	224
19.5.4 处理指令	225
19.5.5 多种语言的交织	225
19.5.6 多语种超级文档	226
19.5.7 多语种 World Wide Web	227
19.6 TEI 编写系统声明	227
19.7 小结	228

第 20 章 编码字符集	229
20.1 集合的乐趣	229
20.2 电话编码：5 位集	229
20.3 ASCII, EBCDIC 和 ISO 646: 7 位集	229
20.4 ISO 8859, ISCII, JIS X 201: 8 位集	231
20.4.1 ISCII	233
20.4.2 JIS X 0201-1979	233
20.4.3 GB 8045-87	233
20.4.4 Adobe 标准编码	233
20.5 扩展 8 位集	233
20.6 16 位集	234
20.7 扩展 16 位集	234
20.8 通用集合	234
20.9 文字	236
20.10 字符集和编码	237
20.10.1 WG4 字符编码模型	237
20.10.2 如何指定字符编码	238
20.10.3 实体	238
20.10.4 元素	238
20.10.5 处理指令	239
20.10.6 存储对象	240
20.10.7 命名字符集和编码	240
20.11 小结	241
第 21 章 断词和分行	242
21.1 空白、单词、连字符和行	242
21.2 断词	243
21.2.1 联接	245
21.2.2 拆分单词并用连字符连接	246
21.2.3 查找	248
21.3 空白	249
21.4 汉语中的断词	251
21.5 小结	252
第 22 章 特殊字符与 SDATA	253
22.1 使用 SDATA 实体	253
22.2 字符的质量保证	256
22.3 重音	256

22.4	HTML 实体	258
22.5	数学脚本和符号	260
22.6	XML	261
22.7	小结	262
第 23 章 从字符到字形		263
23.1	字形映射	263
23.2	使用元素的字形选择	265
23.3	使用实体的字形选择	266
23.4	尺寸	267
23.5	上标和下标	268
23.6	颜色代码	269
23.6.1	黑色、灰色和白色	270
23.6.2	颜色	271
23.7	排版的修饰	282
23.8	小结	282
第 24 章 东亚问题		283
24.1	自定义符号	283
24.2	额外的字符	284
24.3	外来字符与用户定义的字符	285
24.4	定制字体	286
24.5	标记手写文本	287
24.6	“红宝石”注解	288
24.7	本地语言标记	290
24.8	小结和尾注	291

第四部分 附录

附录 A ISO 特殊字符		295
A.1	常用字符	296
A.2	语言	296
A.2.1	ISO 8879:1986//ENTITIES Added Latin 1//EN	296
A.2.2	ISO 8879:1986//ENTITIES Added Latin 2//EN	298
A.2.3	ISO 8879:1986//ENTITIES Greek Letters//EN	302
A.2.4	ISO 8879:1986//ENTITIES Monotoniko Greek//EN	303
A.2.5	ISO 8879:1986//ENTITIES Russian Cyrillic//EN	304
A.2.6	ISO 8879:1986//ENTITIES Non-Russian Cyrillic//EN	307

A.3 符号	308
A.3.1 ISO 8879:1986//ENTITIES Numeric and Special Graphic//EN	308
A.3.2 ISO 8879:1986//ENTITIES Publishing//EN	310
A.3.3 ISO 8879:1986//ENTITIES Diacritical Marks//EN	313
A.3.4 ISO 8879:1986//ENTITIES General Technical//EN	314
A.3.5 ISO 8879:1986//ENTITIES Box and Line Drawing//EN	316
A.3.6 ISO 9573-13:1991//ENTITIES Chemistry//EN	317
A.3.7 ISO 8879:1986//ENTITIES Added Math Symbols:Arrow Relations//EN	319
A.3.8 ISO 8879:1986//ENTITIES Added Math Symbols:Binary Operators//EN	321
A.3.9 ISO 8879:1986//ENTITIES Added Math Symbols:Delimiters//EN	322
A.3.10 ISO 8879:1986//ENTITIES Added Math Symbols:Negated Relations//EN	323
A.3.11 ISO 8879:1986//ENTITIES Added Math Symbols:Ordinary//EN	325
A.3.12 ISO 8879:1986//ENTITIES Added Math Symbols:Relations//EN	326
A.3.13 ISO 8879:1986//ENTITIES Greek Symbols//EN	328
A.3.14 ISO 8879:1986//ENTITIES Alternative Greek Symbols//EN	330
附录 B HTML 特殊字符	332
B.1 HTML 特殊字符	332
B.2 HTML 完全拉丁语 1	334
B.3 HTML 符号	337
B.4 WWW 的草案图标	345
附录 C TEI 特殊字符	347
C.1 -//TEI TRI W4:1992//ENTITIES Basic Arabic Letters//EN	347
C.2 -//TEI TRI W4:1992//ENTITIES Extra Classical Greek Letters//EN	349
C.3 -//TEI TR1 W4:1992//ENTITIES IPA symbols for interchange//EN	355
C.4 -//TEI TR1 W4:1992//ENTITIES Coptic Letters//EN	365
附录 D XML 特殊字符索引	367
D.1 阿拉伯字母	368
D.2 其他阿拉伯字符	372
D.3 阿拉伯数字	373
D.4 亚美尼亚大写字母	374
D.5 亚美尼亚小写字母	375
D.6 其他亚美尼亚字符	376
D.7 孟加拉字母	376
D.8 孟加拉语元音符号	378
D.9 其他孟加拉语字符	378
D.10 孟加拉语数字	379