

Building Scalable and High-Performance Java  
Web Applications Using J2EE Technology

# J2EE Web 应用

## 高级编程

Greg Barish 著 林琪 英宇 译



清华大学出版社  
<http://www.tup.com.cn>  
<http://www.tup.tsinghua.edu.cn>



# J2EE Web 应用高级编程

Greg Barish 著

林 琦 英 宇 译

清华 大学 出 版 社

(京) 新登字 158 号

北京市版权局著作权合同登记号：01-2002-3007

### 内 容 简 介

本书介绍的是用 J2EE 技术构建可扩展和高性能的 Java Web 应用的知识,重点阐述了用 J2EE 规范构建高效的端到端解决方案。本书首先介绍了与可扩展和高性能的 Web 应用有关的基础性知识;然后讲解了用 J2EE 构建 Web 应用的方法和技术,涵盖了客户/服务器通信、客户请求处理、应用服务器的构建和企业应用集成以及数据库设计和操作等知识,对于涉及到的每一种相关技术,均从其可扩展性和性能角度进行了总结和分析;最后本书还对 Web 服务作了总体性的介绍。

本书内容浅显易懂、示例生动、代码丰富,适合于 Web 开发人员和 Java 程序员阅读。

Simplified Chinese edition copyright © 2002 by Pearson Education NORTH ASIA LIMITED and Tsinghua University Press.

Building Scalable and High-Performance Java Web Applications Using J2EE Technology: first publication by Greg Barish, Copyright © 2002.

All Rights Reserved.

Published by arrangement with Pearson Education, Inc., publishing as PH PTR.

This edition is authorized for sale only in the People's Republic of China (excluding the Special Administrative Region of Hong Kong and Macau).

本书中文简体字版由美国培生教育出版集团授权清华大学出版社出版。未经出版者书面许可,不得以任何方式复制或抄袭本书内容。

**版权所有, 翻印必究。**

**本书封面贴有 Pearson Education 出版集团激光防伪标签, 无标签者不得销售。**

**图书在版编目(CIP)数据**

J2EE Web 应用高级编程/(美)巴里希著;林琪,英宇译.一北京: 清华大学出版社, 2002

书名原文: Building Scalable and High-Performance Java Web Applications Using J2EE Technology

ISBN 7-302-05887-3

I.J... II. ①巴...②林...③英...III. JAVA 语言程序设计 IV. TP312

中国版本图书馆 CIP 数据核字(2002)第 071698 号

**出 版 者:** 清华大学出版社(北京清华大学学研大厦, 邮政编码: 100084)

<http://www.tup.com.cn> <http://www.tup.tsinghua.edu.cn>

**责 任 编 辑:** 陈宗斌

**印 刷 者:** 北京昌平环球印刷厂

**发 行 者:** 新华书店总店北京发行所

**开 本:** 787×1092 1/16 **印 张:** 18.5 **字 数:** 473 千字

**版 次:** 2002 年 10 月第 1 版 2002 年 10 月第 1 次印刷

**书 号:** ISBN 7-302-05887-3/TP · 3493

**印 数:** 0001~4000

**定 价:** 38.00 元

# 前　　言

构建 Web 应用时，有两个最基本的问题。首先是开发应用的选择实在太多，“条条大路通罗马”。目前，可以在数以百计的语言和技术中任意挑选，它们都声称能够使 Web 应用的开发更为容易。另一个问题是：如何在部署应用时使之具有高性能和可扩展性没有达成共识，而这却是任何可通过 Internet 访问的 Web 应用所必须具有的共性。

幸运的是，现在已经有了一些统一应用的基础架构，它们不仅可以简化 Web 应用的开发，而且还能支持有效的部署。其中最为常用的方法之一当属 Java 2 企业版(Java 2 Enterprise Edition, J2EE)规范。J2EE 规范分析并扩展了 Java 技术现有的基础，从而使技术人员可以构建支持高性能和可扩展的 Web 应用。J2EE 最为诱人之处在于它对于事务管理、命名和安全等主要的低级平台服务提供了内置的支持，而以往这些服务通常都需要花费大量的设计时间，而且还需要专门的技术方能实现。

尽管 J2EE 将大多数 Web 应用开发的部署进行了统一，但是理解起来却是一件会让人头疼的事情。为了理解 J2EE，要么需要阅读规范本身，要么必须博览众多关于各种特定 J2EE 技术的书(如 Servlet、EJB 和 JDBC 等)。规范本身相当枯燥，利用这一资源往往不具备实践性。如果要把关于每一种技术的所有书都读完，所带来的问题就是花去的时间将会长得使人难以接受，要知道，一些仅仅是介绍 EJB 的书甚至就超过了 1000 页。而且这些书中往往充斥着一些关系不甚紧密的细节问题，还常常会有一些不需要的附录，或者是附录中的内容可以在网上轻松地找到。最后，这些书中一般都没有提及 Web 应用设计的一些不太受关注的领域，例如，数据库设计和网络效率等。

## 本书目标

对于如何利用 J2EE 技术构建高性能和可扩展的 Web 应用，需要有一个实用的概括性总结，而编写本书的目的就是要满足这一需求。在此需要达到一种巧妙的平衡，即一方面要介绍一系列全新的技术以及它们之间的相互关系，还要提供足够的示例以使读者确实了解到 J2EE 究竟如何工作。另一方面又不能过于宽泛，以至于变成一本既没有重点而又厚重的大书(比如，重得甚至于有人会因为举这本书而受伤)。因此，我们的目标很明确，就是要提供一个言简意赅而又实用的概要性介绍。

书中将覆盖 J2EE 的所有主要元素——其规范本身、Servlet、JSP、EJB、消息、JDBC，还有其他一些内容。全书中会提供大量的示例。示例表述得相当简洁，这一点是故意所为，我们的想法是如此精简可以使读者更好地理解如何使用一种技术，而不会因为过于深入到繁杂的细节中去从而把读者搞得头昏眼花。应该知道，关于这些详细的内容，大量的书中都已经有所介绍。如果对 J2EE 中某一部分特别感兴趣，可以参考这些书作为本书的补充材料。

除了要提供一个面面俱到的概括之外，本书还有一个目标，就是想填补有关 Web 应用设计



的空白，它在 J2EE 规范里是没有提及的。尽管规范中说明了如何将不同 Java 技术和 API 加以连接从而构造一个企业应用基础架构，但是并没有谈到相关的一些问题，如网络和数据库设计等。例如，虽然规范中描述了 HTTP 可用于与一个 J2EE 系统通信，但是关于 HTTP 和关系数据库的细节却只字未提。作为一个经验丰富的 Web 应用设计者来说，应该知道设计系统时若要保证高性能和可扩展性，这两者都是至关重要的。

总而言之，本书有如下目标：

- 定义并指出构建可扩展的和高性能 Web 应用时相关的诸多问题。
- 提供一个 J2EE 技术“导航图”，以便按图索骥来设计 Web 应用。
- 简要地描述主要的 J2EE 技术，并重点强调与高性能和可扩展性相关的细节。
- 填补 J2EE 规范所遗漏的 Web 应用设计的空白，如与 HTTP 和数据库设计相关的重要细节内容，这也是与 J2EE 技术有关的最为常用的两个方面。
- 指出不同特定设计方法的优点，并利用实际的性能图表来描述它们之间存在的差异。

最后一点是为了使本书内容更加受到关注。例如，如果说“连接池很好”，这种说法当然也不错，而这也正是很多书所惯用的描述方式。但是如果能够为这一论述提供实际的性能图表作为“证据”，那么将会更具有说服力，同时也更为清晰。本书就是要达到这个目的。

## 本书面向的读者

对于熟悉 Java 的技术人员和设计人员，如果希望构建基于 Java 的 Web 应用从而得到高性能和可扩展性，但同时又对 J2EE 如何用于这些目标以及其底层技术的工作原理不甚了解，那么本书正是为此而编写的。

对于希望了解 J2EE 规范未涉及到的一些内容的读者，这本书同样适用，在此特别强调了目前已经存在的、高效的网络和数据库设计等问题，还讨论了与 XML 和 SOAP 等 Web 服务技术相关的问题，它们在将来可能会出现。

最后，这本书还面向已经对部分 J2EE 技术(如 Java Servlet)比较熟悉，但不了解其他技术(如 Java 消息服务)的读者。

## 有关性能测量的说明

本书中贯穿着多种性能测量和比较的描述。尽管对于每一种体系结构来说，几乎都有着同样的一般趋势(这是因为有关性能趋势的描述是与体系结构无关的)，不过在此把用于测试系统的详细情况列出还是很有帮助的。

所有测试都是在一台带有单个 CPU833 MHz Pentium III，并有 256KB RAM 的 Dell Latitude 计算机上完成的。操作系统和应用软件包括：

- Windows 2000, Professional 版本。
- Apache Web 服务器，1.3.14 版本。
- Java 运行时环境和开发工具包(Java Runtime Environment and Development Kit)，1.3 版本。

- J2EE SDK 及参考实现, 1.3 版本(Beta 版)。
- Apache Jakarta/Tomcat Servlet 容器(Apache Jakarta/Tomcat Servlet Container), 3.2.1 版本。
- Oracle 数据库系统, 8.1.6 版本。

## 开始吧

我希望这本书能够在进行 J2EE 应用开发时提供一个全程的参考。对于如何向您的客户提供应用功能(这其中既包括个人也包括其他企业), 可能需要不断地做出选择。您也许会提出这样的问题, 如: “业务逻辑是放在数据库服务器上好呢, 还是放在应用服务器上? ”, 或者是“我们的批数据需要通过一个 Web 服务器传输呢, 还是通过一个消息服务器传输? ”。暂不考虑与这些选择有关的详细内容, 您肯定希望以某种方式利用这些特性从而提高性能和可扩展性。本书就是要帮助您理解如何在 Web 应用设计中做出权衡, 这包括一般性的和专用的方法, 从而有助于您达到目标。

# 目 录

<b>第1章 可扩展和高性能 Web 应用</b>	1
1.1 Web 应用的出现	1
1.1.1 基本定义	1
1.1.2 Web 的本质特性及其挑战	3
1.2 性能和可扩展性	4
1.2.1 性能	4
1.2.2 可扩展性	6
1.3 Internet 媒体	6
1.3.1 更广泛的受众群体	7
1.3.2 交互性	7
1.3.3 动态性	8
1.3.4 总呈“开放”状态	8
1.3.5 集成性	9
1.3.6 缺乏完全控制	9
1.4 测量性能和可扩展性	10
1.4.1 测量性能	10
1.4.2 测量可扩展性	12
1.4.3 吞吐量和价格/性能比	15
1.5 可扩展性和性能提示	16
1.5.1 考虑端到端	16
1.5.2 可扩展性不等于性能	16
1.5.3 通过比较测量可扩展性	16
1.6 小结	17
<b>第2章 Web 应用体系结构</b>	18
2.1 Web 应用术语	18
2.2 应用需求	19
2.2.1 业务逻辑	19
2.2.2 数据管理	20
2.2.3 接口	20
2.3 Web 需求	21
2.4 抽象 Web 应用体系结构	22
2.4.1 从客户到服务器：瘦客户和胖客户	22



2.4.2 持久性数据管理 .....	24
<b>2.5 N 层应用体系结构 .....</b>	<b>24</b>
2.5.1 客户 .....	24
2.5.2 网络 .....	25
2.5.3 服务器 .....	28
2.5.4 基于层的设计 .....	29
2.5.5 多线程的应用服务器 .....	31
2.5.6 有效中间件带来的问题 .....	32
<b>2.6 可扩展性和性能提示 .....</b>	<b>33</b>
2.6.1 不要对瘦客户期望过高 .....	33
2.6.2 使用或建立多线程应用服务器 .....	34
2.6.3 确定合适的粒度 .....	35
<b>2.7 小结 .....</b>	<b>35</b>
<b>第 3 章 J2EE 规范 .....</b>	<b>37</b>
3.1 规范概述 .....	37
3.2 部署问题 .....	40
3.2.1 包装 .....	40
3.2.2 部署描述符文件 .....	41
3.3 平台技术与服务 .....	43
3.3.1 通过 RMI-IIOP 实现组件通信 .....	43
3.3.2 使用 Java 事务 API 实现事务管理 .....	45
3.3.3 实现资源查找的 JNDI .....	46
3.4 J2EE 和体系结构 .....	48
3.5 小结 .....	49
<b>第 4 章 可扩展性和性能技术 .....</b>	<b>50</b>
4.1 缓存与复制 .....	50
4.2 并行 .....	55
4.3 冗余 .....	58
4.4 异步 .....	59
4.5 资源池 .....	61
4.6 小结 .....	67
<b>第 5 章 HTTP 客户/服务器通信 .....</b>	<b>69</b>
5.1 HTTP 协议 .....	69
5.2 部署模式 .....	71
5.2.1 带有浏览器客户的应用 .....	71
5.2.2 不带浏览器的应用 .....	71

5.3 HTTP 效率 .....	72
5.4 HTTP 详细内容 .....	73
5.4.1 语义 .....	73
5.4.2 HTTP 请求 .....	75
5.4.3 GET 方法 .....	75
5.4.4 POST 方法 .....	79
5.4.5 HTTP 1.1 缓存 .....	80
5.4.6 连接管理 .....	83
5.5 可扩展性和性能提示 .....	85
5.5.1 理智地使用 GET 和 POST .....	85
5.5.2 对于非浏览器客户考虑 HTTP .....	85
5.5.3 提升 HTTP 响应缓存 .....	85
5.5.4 支持持续连接 .....	87
5.6 小结 .....	87
<b>第 6 章 请求处理 .....</b>	<b>88</b>
6.1 一般问题 .....	89
6.2 特定问题 .....	89
6.2.1 连接管理 .....	91
6.2.2 数据编组 .....	91
6.2.3 请求服务 .....	92
6.2.4 缓存环境中的数据本地性 .....	94
6.3 请求处理模式 .....	94
6.3.1 同步通信 .....	95
6.3.2 异步通信 .....	95
6.3.3 可扩展性和性能问题 .....	98
6.4 请求处理和 J2EE .....	99
6.4.1 Web 服务 .....	99
6.4.2 利用 Java servlet 和 JSP 实现同步处理 .....	100
6.4.3 使用 Java 消息服务实现异步处理 .....	100
6.5 可扩展性和性能提示 .....	101
6.5.1 建立异步解决方案 .....	101
6.5.2 线程间的流数据 .....	102
6.5.3 开发有效的远程接口 .....	107
6.6 小结 .....	110
<b>第 7 章 基于 Java servlet 的会话管理 .....</b>	<b>111</b>
7.1 生成动态响应 .....	111
7.1.1 公共网关接口 .....	111



7.1.2 通过 API 扩展 Web 服务器 .....	112
7.1.3 重定向 Web 服务器请求 .....	112
7.2 使用 servlet .....	113
7.2.1 servlet 和 servlet 容器 .....	113
7.2.2 与 servlet 交互 .....	114
7.2.3 Web 服务器与 servlet 容器集成 .....	116
7.3 开发 servlet .....	117
7.3.1 设计 servlet 接口 .....	117
7.3.2 建立 servlet 的代码 .....	118
7.4 servlet 执行 .....	121
7.4.1 servlet 容器 .....	121
7.4.2 servlet 和多线程 .....	122
7.5 servlet 和会话管理 .....	125
7.6 部署 servlet .....	130
7.7 使用 JSP 开发 servlet .....	132
7.7.1 JSP 页面示例 .....	132
7.7.2 JSP 页面的结构 .....	133
7.7.3 JSP 如何工作 .....	135
7.7.4 JSP 指示 .....	138
7.7.5 JSP 到底是什么 .....	139
7.8 可扩展性和性能提示 .....	139
7.8.1 使用细粒度的串行化 .....	139
7.8.2 使用基于硬件的负载平衡 .....	140
7.8.3 使用 servlet 实现会话管理，而非业务逻辑 .....	140
7.8.4 再三考虑 JSP .....	141
7.9 小结 .....	141
<b>第 8 章 利用企业 JavaBean 构建应用服务器 .....</b>	<b>143</b>
8.1 应用服务器的需求 .....	143
8.2 企业级 JavaBean：J2EE 解决方案 .....	144
8.3 EJB 的工作原理 .....	145
8.4 EJB 类型 .....	146
8.5 应用示例 .....	147
8.6 EJB 设计 .....	148
8.6.1 会话 bean .....	148
8.6.2 实体 bean .....	148
8.6.3 消息驱动 bean .....	153
8.7 EJB 实现 .....	153

8.7.1 会话 bean.....	153
8.7.2 实体 bean.....	160
8.7.3 实体 bean 和 EJB 2.0.....	161
8.7.4 消息驱动 bean.....	172
8.8 客户/EJB 集成 .....	175
8.9 可扩展性和性能提示 .....	179
8.9.1 尽量用消息驱动 bean 而不是会话 bean .....	179
8.9.2 使用无状态会话 bean.....	180
8.9.3 尽量采用粗粒度的 EJB 方法 .....	181
8.9.4 要么很好地使用 BMP，要么干脆不用 .....	183
8.9.5 了解您的开发商.....	184
8.10 小结.....	184
<b>第 9 章 基于消息实现高效的企业应用集成 .....</b>	<b>186</b>
9.1 B2B 型的工作实例.....	186
9.2 Java 消息服务.....	187
9.3 JMS 概念.....	187
9.3.1 提供者 .....	188
9.3.2 客户.....	188
9.3.3 消息.....	188
9.3.4 管理对象.....	190
9.4 JMS 编程模型.....	190
9.4.1 特定于模型的管理对象接口 .....	191
9.4.2 消息使用的同步性 .....	191
9.5 JMS 可靠性与性能 .....	192
9.5.1 客户确认 .....	192
9.5.2 消息持久保存 .....	193
9.5.3 时间依赖性和 JMS 发布模型 .....	193
9.6 一个 JMS pub/sub 应用示例 .....	194
9.6.1 开发消息发布者 .....	194
9.6.2 开发消息预约者 .....	196
9.6.3 关于部署 .....	199
9.7 可扩展性和性能提示 .....	199
9.7.1 使用消息 .....	199
9.7.2 理解 JMS 效率-可靠性的折衷 .....	202
9.8 小结 .....	202
<b>第 10 章 高效的数据库设计 .....</b>	<b>204</b>
10.1 数据库技术和关系模型 .....	205



10.2 逻辑数据库设计 .....	206
10.3 物理数据库设计 .....	207
10.3.1 表和行 .....	208
10.3.2 约束 .....	209
10.4 查询数据库 .....	209
10.4.1 查询数据 .....	209
10.4.2 嵌套查询 .....	211
10.4.3 连接查询 .....	211
10.5 其他重要的数据库对象 .....	213
10.5.1 视图 .....	213
10.5.2 存储过程 .....	214
10.5.3 触发器 .....	216
10.5.4 索引 .....	217
10.5.5 序列 .....	218
10.5.6 其他对象 .....	219
10.6 查询处理 .....	219
10.7 可扩展性和性能提示 .....	222
10.7.1 理解如何使用数据库 .....	222
10.7.2 理解何时使用数据库 .....	223
10.7.3 理解如何访问数据 .....	224
10.7.4 规范数据模型 .....	225
10.7.5 有选择地实现模型的非规范化 .....	229
10.7.6 使用存储过程 .....	231
10.7.7 避免触发器及其他隐式执行 .....	234
10.7.8 了解开发商 .....	234
10.8 小结 .....	235
<b>第 11 章 使用 JDBC 和 SQL 高效查询数据库 .....</b>	<b>236</b>
11.1 使用 JDBC 的原因 .....	236
11.2 JDBC 概念和对象 .....	236
11.2.1 相关 JDBC 对象及其关系 .....	237
11.2.2 连接数据库 .....	238
11.3 编写 JDBC 查询 .....	238
11.3.1 处理语句 .....	238
11.3.2 循环处理结果 .....	239
11.3.3 执行单个更新 .....	242
11.3.4 其他类型的更新：创建表和存储过程 .....	242

11.4	更高级的问题 .....	243
11.4.1	准备语句(prepared statement) .....	243
11.4.2	动态 SQL .....	244
11.4.3	事务管理 .....	245
11.4.4	双向结果循环 .....	247
11.4.5	可更新结果 .....	248
11.4.6	执行批更新 .....	248
11.5	可扩展性和性能提示 .....	249
11.5.1	在可能的情况下使用 PreparedStatement .....	249
11.5.2	对一个远程数据库使用批更新 .....	251
11.5.3	不要过分使用提交 .....	252
11.5.4	使用多线程实现并行查询 .....	253
11.6	小结 .....	253
<b>第 12 章 Web 服务: Web 应用的未来 .....</b>		<b>254</b>
12.1	Web 服务的实际使用 .....	255
12.2	Web 服务到底是什么 .....	255
12.3	Web 服务技术 .....	257
12.3.1	概述 .....	257
12.3.2	综合 .....	258
12.4	XML: 自描述数据 .....	259
12.4.1	DTD 和模式语言 .....	260
12.4.2	解析 XML .....	261
12.4.3	与 XML 相关的技术 .....	266
12.5	开发 Web 服务 .....	266
12.6	使用 WSDL 描述 Web 服务 .....	267
12.6.1	定义 .....	267
12.6.2	示例 .....	268
12.7	使用 SOAP 调用 Web 服务 .....	270
12.7.1	SOAP 如何工作 .....	270
12.7.2	使用 HTTP 上的 SOAP .....	271
12.8	利用 UDDI 注册 Web 服务 .....	273
12.8.1	标准 .....	274
12.8.2	UDDI API .....	274
12.9	重览全局 .....	276
12.9.1	提供者角度 .....	276
12.9.2	使用者角度 .....	277



12.10 可扩展性和性能问题 .....	277
12.10.1 远程方法的复制与负载平衡 .....	277
12.10.2 XML 解析性能 .....	277
12.10.3 解析与查询 XML .....	278
12.11 小结 .....	280

# 第1章 可扩展和高性能Web应用

## 1.1 Web 应用的出现

仅仅是在几年之间，在世界范围内，无论是信息的提供方式还是使用方式都因 Internet 而发生了改变。其硬件和软件技术使得每个人不仅能够成为信息的使用者，而且几乎所有人都能够作为信息的提供者。Internet——特别是 World Wide Web(Web)——在非常短的时间内就已经被公众认为是重要的信息共享的平台，许多组织都在尽力创建有用的 Web 应用，从而为使用者提供更大的价值。

这些 Web 应用允许使用者在线地购买书和光盘。它们使得企业可以使用 Internet 来实施安全的事务处理。工人利用 Web 应用寻找工作；老板利用 Web 应用寻找雇员；使用由经纪人提供的在线应用可以购进和抛售股票；旅行者则可以利用 Web 应用来预订机票和宾馆。这样的例子还有很多很多。很明显，目前无论是在公共的 Internet 上，还是在数不胜数的公司内部网 (Intranet)中，都存在着大量有用的 Web 应用。

本书介绍构建具有高性能和可扩展性的企业 Web 应用所需的通用技术。一般来说，这表示所构造的应用不仅要有合理、稳定快速的运行速度，而且对于不断增长的用户和请求需求还应具有较强而渐增的容许性。尽管我们会花很多时间来考虑这个一般性的问题，但是我们要讨论的核心却是围绕着 Java 2 企业版(Java 2 Enterprise Edition, J2EE)规范来建立一个解决方案。在对建立这一类应用的具体细节深入分析之前，就目前而言，很重要的一点是需要明确和理解整个问题。更准确地说，就是要定义 Web 应用和可扩展性，这是非常重要的。

### 1.1.1 基本定义

本书中，Web 应用有一个非常通用的定义——这是一种通过 Internet 技术加以连接的客户/服务器软件，可以传输其处理的数据。通过“Internet 技术”，我所指的是在信息的使用者和提供者之间，组成相应网络基础架构的硬件和软件的集合。Web 应用可以通过专门的客户端软件来访问，也可能利用一个或多个有关的 Web 页面访问，这些页面基于某种特定的用途可进行逻辑分组。这里所说的用途可以是任意一件事情，例如，可以是买书、处理股票订单，也可能仅仅是作为让使用者阅读的内容。

注意，我们所讨论的是 Web 应用，而不只是“Web 网站”。实际上，这二者之间的区别对于理解本书的一个关键主题相当重要。大多数非技术人员往往不会区别 Web 网站和 Web 应用。除了术语的提法不同之外，对他们而言，其作用都是一样的，都可以使之实现在线购书、在线预订机票、在线购票等操作。

不过，如果您是一名技术人员，那么这两者之间还是存在差别的。对您来说，如果有人谈



到类似于一个 Web 网站的性能问题时，您可能会开始想到其后台的具体细节。我也不例外。您会考虑所运行的是 Apache 还是 IIS，还有它使用的是 Java servlet、PHP 还是 CGI-bin Perl 脚本。技术人员和非技术人员在思路上所存在的不同可能会造成一定的误解。技术人员，往往习惯性地把“Web 网站”与服务器端联系起来。而与此同时，我们都知道，Web 应用并不仅仅包括服务器端；它还需要有网络和客户端。因此，从这个意义上说，一个 Web 网站(服务器)与 Web 应用(客户机、网络和服务器)并不是一回事。

虽然这本书所强调的是服务器端解决方案，在此还会讨论客户端和联网的有关内容，这是因为对于终端用户如何理解 Web 应用，它们有着很重要的影响。也就是说，我们会讨论 Web 网站内端对端的交互，它意味着由客户到服务器，再返回到客户。这自然是一个重点。毕竟，大多数使用 Web 的人所关心的就是它的端到端的操作。如果在网上购买音乐会的门票需要花费他们较长的时间，则造成这一现象的原因可能有多种，可能是因为一个速度较慢的调制解调器(MODEM)所致，或者原因在于一个负载过重的服务器，也许还可能是由于网络阻塞所造成。无论原因如何，结果都是一样的，即“慢吞吞”的应用会耗费大量的时间。作为技术人员，我们所关注的将不仅仅是对于一个用户应用可能比较慢，更重要的则是随着访问量的增加，系统也许会越来越慢。

既然我们对于 Web 应用的范围已经有了更为明确的认识，下面再来看它的核心组件。任何一个在线应用都有一些主要的组成部件，每个部件即代表着一种可能性——也许是一个问题，也许是一个挑战，这要取决于您如何看待它。可能您对这些组件已经比较熟悉了，不过在此确保每个人都明白是一件有益而无害的事情，特别是这些术语将贯穿本书始终。下面就先从客户端开始见图 1-1。

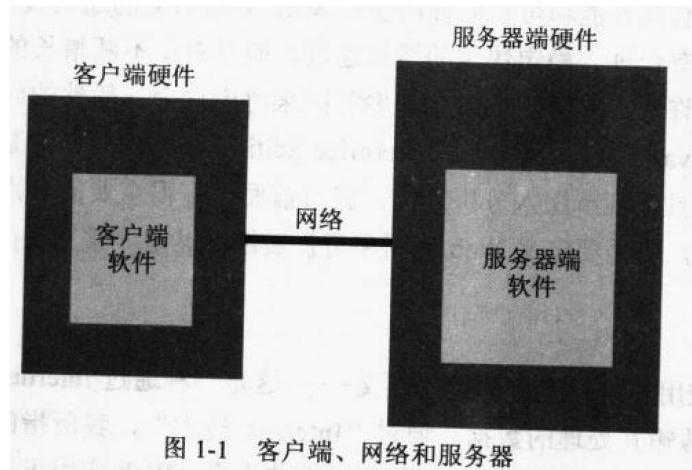


图 1-1 客户端、网络和服务器

我们可以说使用者是通过客户端软件(即使用 Web 来检索和处理数据的 Web 浏览器或应用)来使用 Web 应用的，这些客户端软件运行在客户端硬件(即 PC、PDA 等等)之上，而应用数据的提供以及生产者对具体处理的控制则均通过服务器端软件(即 Web 服务器、服务器端组件软件、数据库等)来实现，这些服务器端软件运行于服务器端硬件(即高端多处理器系统、集群等等)。将客户端与服务器进行连接(从调制解调器或客户端设备的网络端口到服务器端的网络设备)就形成联网基础架构。图 1-1 用图的形式描述了客户端/服务器的关系。注意服务器端要大一些，通常情况下，我们假设服务器端有更多可支配的资源。

在此，需要强调的很重要的一点，即作为服务器软件的一种，要把 Web 服务器特别区别出

来，因为它在调度客户端与服务器之间的通信(HTTP流量)时总是起着一个核心的作用。在本书中，当提到“服务器端”时，一般都包括Web服务器。如果有必要将它与服务器端的其他软件相区别时，我们还会专门指出。

### 1.1.2 Web的本质特性及其挑战

尽管由于Web应用的存在使得Internet迅速成为一个应用广泛的媒体，不过，不理解Internet的本质特性，也会造成一些技术上的难题。即使是最基本的问题，例如，使信息的提供者快速而可靠地为所有需要的人发送相应的信息，这个问题就既不简单也不好理解。与其他挑战一样，这个问题的复杂性在于必须处理Internet媒体的本质特性。基于许多原因，Internet与诸如广播、电视和报纸等其他信息共享模式存在着显著的区别。这其中最为重要的两个原因可能是，一方面它所面向的群体实在过于广泛(客户的数量不可估计)，而另一方面，所面向的客户群体在任何时刻都可能对任意的提供者提出信息请求(工作需求量不可估计)。

与其他媒体不同，Internet信息提供者不具有预先了解其客户的能力。例如，报纸能够在印刷每一期之前知道它的发行量。另外，这些媒体还有一个优势，就是可以控制其增长速度，以确保每天有足够的员工来发放报纸，而且在最后期限之前的头一天晚上可以有足够的资源和时间来准备第二天早上的报纸发放工作。除此之外，报纸也不会遇到发行量突然跳跃式增长的情况。与Internet相比，即使是大城市的大型报纸，其增长也不过是渐进式的。例如，《华盛顿邮报》始创于1877年，当时其发行量为10000。到了1998年，每天的发行量已经接近了800000份，而且周日版的发行量还要多。<sup>\*</sup>其平均增长速度为每年不少于6500份，即每天17份。

Web应用的开发人员对于Web应用的增长速度却是既爱又恨。一方面，如果每天都有17位新的用户，他们当然希望这样的增长速度。要是能够为增长过快而考虑适当扩展，该是多么幸福！这样您就不用每天工作到晚上9:30，而是在下午5点就轻松地回到家里。同时，如此高的增长率也正是人们对Web应用这样着迷的原因之一，因为利用它只在几秒钟之内就可以达到世界上的几乎每一个角落。Web应用的增长率可能达到数百万之多。尽管在企业方面这是一个很好的预兆，但同时也带来一个严重的问题，该如何处理这么多的需求呢？

在Internet上，发行量就相当于“页面的点击率”，也就是说，对于某个文档的请求数目。页面点击率可能会在一夜之间出现“疯长”。Starr报告的在线公布应该算是Web领域中对此最适当的一个例子。因为大多数美国人都知道，这个报告是克林顿在任期间由独立顾问处所整理的。我们可以毫不夸张地说，无论是美国公众还是国际报业集团都迫不及待地想对其有所了解。

Starr报告于1998年夏天在政府Web网站上首次在线发布时，数以万计的人都试图进行下载。作为Internet主要ISP之一的Sprint公司，其代表报告称带宽需求高峰值较之平常要高出10到20个百分点；AOL的一个代表则表示“突然出现了30个百分点的增长”；NetRatings，这是一个像Nielsen一样的Internet上知名度很高的公司，据它估计，某些时候，每5个Web用户中就有一个在请求此报告，或者是相关的新闻。关于这一事件以及1998年秋天有关Internet可扩展性的林林总总，CNET.COM上有大量相应的报道。

<sup>\*</sup> 资料来源：<http://www.thewashingtonpost.com>