

9.6.1 大格局

Web 发展的初期，搜索引擎就是一些相当简单的数据库，可以帮助用户在 Web 上定位文档。现在，Web 上有数十亿可供访问的页面，搜索引擎已经成为因特网用户查找信息不可缺少的工具。它们在不断地发展，以应对 Web 庞大的规模，因此，现在已经变得相当复杂了。

面对数十亿的 Web 页面，和数百万要查找信息的用户，搜索引擎要用复杂的爬虫来获取这数十亿 Web 页面，还要使用复杂的查询引擎来处理数百万用户产生的查询负荷。

我们来考虑一下产品级 Web 爬虫的任务，它要获取搜索索引所需的页面，它要发出数十亿条 HTTP 请求。如果每条请求都要花半秒钟的时间（对有些服务器来说可能慢了，对另一些服务器来说可能快了²⁷⁾），（对十亿份文件来说）就要花费：

$$0.5 \text{ 秒} \times (100\,000\,000) / (60 \text{ 秒} / \text{天}) \times (60 \text{ 分} / \text{小时}) \times (24 \text{ 小时} / \text{天})$$

如果请求是连续发出的，结果差不多是 5700 天！很显然，大型爬虫得更聪明一些，要对请求进行并行处理，并使用大量机器来完成这项任务。但由于其规模庞大，爬行整个 Web 仍然是件十分艰巨的任务。

9.6.2 现代搜索引擎结构

现在的搜索引擎都构建了一些名为“全文索引”的复杂本地数据库，装载了全世界的 Web 页面，以及这些页面所包含的内容。这些索引就像 Web 上所有文档的卡片目录一样。

242

搜索引擎爬虫会搜集 Web 页面，把它们带回家，并将其添加到全文索引中去。同时，搜索引擎用户会通过 HotBot (<http://www.hotbot.com>) 或 Google (<http://www.google.com>) 这样的 Web 搜索网关对全文索引进行查询。Web 页面总是在不断地发生变化，而且爬行一大块 Web 要花费很长的时间，所以全文索引充其量也就是 Web 的一个快照。

现代搜索引擎的高层结构如图 9-7 所示。

9.6.3 全文索引

全文索引就是一个数据库，给它一个单词，它可以立即提供包含那个单词的所有文档。创建了索引之后，就不需要对文档自身进行扫描了。

注 27：这取决于服务器的资源、客户端的机器人，以及两者之间的网络状况。