一种分布式搜索引擎设计*)

A Distributed Search Engine Design

印鉴! 邹胜!

(中山大学计算机科学系 广州510275)(深圳证券交易所 深圳518010)2

Abstract This paper presents a distributed search engine design of an on-line bookstore system. Several principles are introduced such as database miniaturization, the entire structure and the main modules are explained in detail. Compared with a centralized structure, the distributed structure has several advantages including high speeds efficient usage of network bandwidth, less security problems, etc. The system uses feedback of the users to judge the information quality, select search engines and update databases. So, the system performances are improved.

Keywords Search engine Distributed structure On-line bookstore system

1 引言

随着 Internet 的发展,地理上分散的、功能独立的计算机系统内的梅量信息资源也由封闭式转变成开放式。其信息特点是:(1)Internet 的信息组织形式各异,分布广泛;(2)数据和服务的类型都在增加,可利用性和可靠性也在不断的变化;(3)信源的动态性。这些特点导致信息量巨大,而且信息的获取并非容易。由此带来了一个重要的问题就是搜索引擎的设计。所谓搜索引擎,简单地说,就是指对 WWW 站点资源和其他网络资源的检索和管理的一类检索系统机制。

搜索技术的发展是伴随计算机应用而生的,从最初的文件检索和文档的查找,到现在的 Internet 搜索。在这过程中,随问题的提出而逐步得到解决,例如早期人们在从大量的文件和大篇幅的文档中找到自己所需的文件和字符串。在 Internet 中人们的视野显得更开阔了,但是人的需求与自身的"能动力"有相当大的差距。这就要求有一种智能化的搜索机制来完成。从20世纪80年代起人们就开发了诸如 Archive、WAIS、Veronica等检索工具,从90年代中期起又出现了检索万维网信息资源的搜索引擎技术,并以此构造检索所有各类网络信息资源的集成化支撑体系。例如 Yobno、Alta Vista、Infoseek、Excite 和搜狐等。

本文首先对搜索引擎作了一个简单的介绍,然后 以我们所设计的一个网上购书系统为例,具体介绍了 其搜索引擎的设计。

2 搜索引擎

2.1 基本框架

搜索引擎的概念模块包括:客户端、数据库、Robot 系统。客户端的设计是一个比较个性化的工作。数据库的关键是客户端和 Robot 系统间的接口问题。作为整个系统而言,数据库的技术是成熟的。整个技术的关键在于 Robot 系统的提高,它是整个搜索的关键环节。图1是一个典型的搜索引擎概念图。

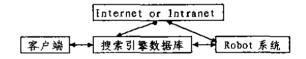


图1 搜索引擎概念图

2.2 搜索引擎的分类

搜索引擎的分类方法有很多种,这里按运行机制 来分类:

- (1)集中式搜索引擎 通常由三部分组成:客户端、数据库、获取网络信息的 Robot。客户端查询数据库,Robot 则帮助扩充和更新数据库。这是一个独立的单结构搜索引擎。
- (2)多体搜索引擎 是一种建立于各种集中式搜索引擎之上的网络信息搜索工具。它根据用户的查询要求,启动多个集中式搜索引擎查询用户需要的信息,并对所有返回结果进行核查、整理、综合,然后返回给用户。它是一个由不同独立分支组成的系统。

^{*)}本文研究得到国家自然科学基金部分资助和广东省自然科学基金资助。

(3)代理(Agent)技索引擎 在收到用户提出的 查询请求后,先在本地数据库内搜索;如果找不到用户 需要的信息,则再向某一个集中式的搜索引擎发出请求,得到查询结果后,将结果在返回给用户的同时存入 本地数据库以备下一次查询。

2.3 搜索引擎的组成

- 1. 客户端 客户端是人机交互的界面部分,从这一部分可以判断出该搜索引擎的功能,除了人的直接 视觉思维受影响外,实质性的搜索过程是不被人所感知的、因此,这一部分的工作除了做好与数据库的接口外,就是界面的布局等工作。搜索引擎提供的信息,不仅要广,还有准确、快捷。
- 2. 技索引擎数据库 数据库的选择是很重要的、它关系到数据的存储、检索、数据的维护等等。如果数据库选择不恰当,将会影响到整个搜索效果。当然要关注数据库与其它两个部分的接口。
- 3. Robot 系统 Robot 系统又叫 Spider,它的主要作用是创建并修改搜索引擎使用的数据库(不是唯一的途径),它获取信息要执行以下三项步骤(或是其中的一到两项);
- (1)探索:通过访问每一个可能的 IP 地址判断 Internet 上的站点。
 - (2)站点爬行:研究某一站点及其内容。
 - (3)编索引:搜集网页上所有的信息。

3 网上购书系统的搜索引擎设计

下面介绍我们所设计的一个网上购书系统的搜索引**整**。

3.1 设计原则

设计思想的核心是:通过合理的结构与运行机制,使系统性能和服务质量达到一定的水平。根据网络信息搜索的特点和网络现状,并充分考虑了搜索引擎的发展趋势,提出了以下设计准则:

(1)分布式整体结构 现有的书籍查询系统大都采取集中式结构,但在运行中逐渐暴露出一些问题。首先集中式的搜索引擎不能合理地利用网络带宽,既要为大量用户提供查询服务,又要随时对大量的数据进行更新,这将形成网络瓶颈。其次所有用户查询一个中心数据库,会给系统带来很大负担。集中式系统还存在安全性、可靠性的问题。

为了解决以上问题,选用了分布式的系统结构,各个节点拥有自己的子系统。每个子系统拥有一个搜索引擎,负责查询本地信息,各搜索引擎可以相互协作,构成分布式的搜索引擎系统。这样的设计具有合理的带宽分配和系统负载分配;各个区间的用户都将具有较快的查询速度;而且系统的安全性和可靠性也将更

有保证。

(2) 數据库小型化 现在大多数著名的搜索引擎系统,大多向着大而全的方向发展,力争将自己的搜索范围扩大到整个 Internet。然而由于数据库规模庞大,为了提高查询速度往往牺牲了信息质量,而且这样的系统需要性能极高的硬件环境。同时对于整个网络的搜索还会伴随其它的阻碍因素,语言系统其实是一个很严重的问题,虽然本系统不存在这样的大问题,但对于设计而言是不可不考虑的。

对于用户而言,最在意的是查询的查准率,而不仅仅是查出率。也就是说,我们并不一定需要得到所有符合查询要求的信息,对于一次查询,大型搜索引擎一般会得到成千上万条符合要求的信息,而人们只会浏览其中的几十条。因此,对于用户来说,能够返回成千上万条信息的搜索引擎与只能返回几百条信息的搜索引擎并没有什么区别,提高信息质量才是关键。

基于以上几点考虑,提出了数据库小型化的原则。每个搜索引擎自带的数据库只存放本地用户经常查询的信息。在满足绝大多数用户查询要求的前提下,限制数据库的规模。这样可以提高查询的速度并且不会使系统负担太重。对于一些本地数据库不能满足的查询要求,系统将求助于其他搜索引擎以获取查询结果。

(3)有选择地协作 各个节点的搜索引擎之间如何协作是系统设计需要解决的重要问题,这个问题主要包括两个具体问题;a. 本地搜索引擎何时需要其他搜索引擎的帮助;b. 需要其他搜索引擎帮助时、向哪些搜索引擎发出帮助请求,很直观地,希望当本地搜索引擎无法满足用户需求时,系统将自动地求助于其他搜索引擎,并且出于对查询速度和效率的考虑,不希望将求助信息广播式地发给所有搜索引擎,而希望系统能够选择最有可能满足用户查询请求的搜索引擎、向它们发出请求。这样才能建立各搜索引擎之间有效的合作。

(4) 充分利用用户反馈 为了有效地提高信息服务质量,利用用户反馈是必要的,数据库小型化和选择协作对象都需要用户反馈信息的帮助。只有知道了用户的需求,才能在数据库规模不大的情况下存贮最有用的信息,才能选择最有可能满足用户查询请求的查询对象。

3.2 系统结构

网上购书系统的搜索引擎同其他的搜索引擎总体结构是一样的,其大致的结构是:用户界面、查询对象选择模块、数据库、远程查询接口和 Robot。

(1)用户界面 作为用户和系统内部结构的接口, 主要负责接收用户的请求、将结果返回用户,以及获取 用户的反馈。这些功能由相应的查询界面、结果综合、 用户行为监控3个子模块完成。

查询界面负责接收用户的各种查询请求,并进行初步的处理,形成系统内部规范的查询条件。

本地搜索引擎的查询结果与远程搜索引擎的查询 结果,都要经过结果综合模块的再加工,去除无效或重 复链接、并依匹配度重新排序、统一表达方式,才能形成最后的结果页面,返回给用户。

用户行为监控模块负责获知用户对系统提供的信息的满意程度,即获取用户的反馈信息,并将获取的有关用户的知识送给有关模块(搜索引擎数据库、查询对象选择模块),使这些模块能对用户需求进行不断的跟踪或者帮助这些模块进行决策,具体说,用户行为监控模块可以利用用户的点击信息,记录用户在查询结果中的选择。用户如果访问了系统提供的某个网络资源(有点击动作发生,或是同一用户的点击次数来计,或是新增用户的速度来计),系统就会记录下来,并且理解为用户对这条信息比较满意,还可以以用户的评价结果来衡量,这就要求在系统中增加分析部件,

- (2) 查询对象选择模块 负责判断对于当前的用户请求,哪些搜索引擎最有可能提供令人满意的查询结果。用户的反馈意见和各个搜索引擎的历史表现会帮助它决定将查询请求发送给本地搜索引擎还是远程搜索引擎。如果是远程搜索引擎,又是哪一些。
- (3)數据库 本系统的各个节点的数据库保存了用户经常查询的信息,可以满足绝大部分的本节点的查询请求,同时可以为远程搜索引擎提供查询服务。本系统中,搜索引擎数据库记录了在每一次查询中每一个搜索引擎的表现值,系统将依据这些数据为新的查询选择适当的搜索引擎。
- (4)远程查询接口 当本地搜索引擎不能满足某些用户查询要求时,需要求助于其他搜索引擎,远程查询接口使系统可以向远程搜索引擎发出查询请求,也可以接收其他搜索引擎发来的请求,并将请求交给本地数据库进行处理。
- (5)Robot Robot 为本地数据库搜集网络资源, 并为信息更新提供帮助,将网络资源的变化及时通知 数据库(本系统的所有数据库数据都是手动获取)。

3.3 系统实现

由于本系统只针对网上购书、因此数据的采集编排较为单一,下面给出几个关键内容的设计:

(1)关键词规范 显然机器的辨识能力不同于人,机器想要很好地完成任务,必须给出相应的规范。本系统定义了三种表,语义编码表、词性表、书籍分类表。这三个表的定义遵从生活习惯,后两者相对稳定,前者的变数较多。

(2)对搜索引擎的合作协议 台理的协作机制是 • 76 • 保证查询质量和查询速度的必要条件,也是分布式系统的设计重点。对于分布式的搜索引擎来说,首先要判断何时需要协作,其次是选择适当的搜索引擎进行协作。本系统的特点是,充分利用用户的反馈,帮助完成协作中的决策,查询对象选择模块是完成搜索引擎之间相互协作的核心部分。与用户行为监控模块及搜索引擎数据库配合,它能确定将查询请求发给本地数据库还是远程搜索引擎。如果要发给远程搜索引擎,它将替用户选择适当的搜索引擎。下面将对这一部分的实现作出具体说明。

已经提到,当本地搜索引擎无法满足用户需要时,将求助于其他搜索引擎。那么,如何判断用户的需求是否被满足了呢?通过对用户行为的监控,以及一些简单而又合理的规则来对用户的满意程度作出判断并采取相应的措施。具体方法是:

在搜索引擎的查询结果中,每一条结果对查询请求有一个内容匹配度,对每一个用户的查询请求,定义一个内容匹配度域值 r.根据查询需求,系统首先在本地数据库中查询,如果本地数据库内有满足此查询请求的信息,即查询结果的内容匹配度大于 r,表示用户的需求很可能在本地数据库内得到满足、则返回从本地数据库中查询的结果。如果数据库内没有满足此查询请求的信息,即查询结果的匹配度小于 r,或用户在向后翻看结果时,所显示的结果的匹配度小于 r 时,系统将启动远程查询接口要求远程搜索引擎帮助查询,共同产生查询结果、反馈给用户。

在未启动远程查询时,用户在前一页查询结果中没有找到想要的信息(用户行为监控模块没有发现用户的点击动作),在用户向后翻看时,系统将这个查询的 r 增加;当用户在前一页查询结果中找到很多的信息(用户行为监控模块发现用户的点击动作很多),系统将这个查询的 r 减小。调节的幅度可以由系统根据实际使用情况确定。

这些规则保证了本地数据库的优先地位,同时,通过对用户行为的监控,在用户对信息质量不满意时切换查询对象,这在兼顾了查询效率的同时,考虑了用户对信息质量的要求,将对提高信息质量有一定的帮助。

(3)对数据库中技索引擎的管理 首先利用搜索引擎数据库来记录以往查询中各个搜索引擎的表现。利用查询关键词为每个搜索引擎作一个索引,具体方法是:在搜索引擎数据库内存贮了一个 t×n 的二维矩阵,t 为关键词表的长度,n 是搜索引擎的数量。矩阵中的每一个元素代表着所在列的搜索引擎对所在行的关键词的查询能力。对于一个含有 m 个关键词的查询请求、如果某个搜索引擎提供的信息被用户选中(用户行为监控模块发现用户的点击动作),则将这个搜索引擎

对应的矩阵列上,这 m 个关键词对应的元素的分值增 加。这样这个矩阵中的每一列就反映了一个搜索引擎 的综合查询能力。系统对增加新的关键词、合并重复的 关键词,进行统一管理,保证各子系统统一,

(4)数据库更新 首先,要保证数据库内信息的有 效性(信息源是否被更新或取消)。这只需要令 Robot 定期将数据库内的信息与信息源进行核对即可。

其次,系统在索引本地资源的同时,对符合用户信 息需求的外部信息进行缓存(即存储用户点击的外部 页面索引),出于对系统性能的考虑,每一个本地数据 库中的信息量都不能超过一定的上限,因此本地搜索 引擎只能保存一部分网络信息,那么最优的选择就是 保留那些本地用户最需要的网络资源、然而很难在系 统刚刚开始使用时就准确地知道本地用户最需要哪些 信息。而且系统在使用过程中,由于网络资源的不断变 化,用户组成的不断变化,用户兴趣的不断变化等等, 用户最需要的信息也会不断变化,因此需要不断更新 数据库内存储的信息,以便跟上用户需求的变化。

利用用户的反馈来跟踪用户需求的变化,并不断 淘汰数据库内过时的信息,添加新的信息。具体的做法 是: 当一条新的信息被存入数据库时,它被赋予一定的 "生命力";当它符合某个用户的查询要求(被点击)时, 便增加其生命力;生命力随时间的流逝而减小,根据系 统容量,系统保留生命力大的记录。这种机制保证了数 据库只存贮用户经常需要的信息。

(上接第73页)

比高的编码算法。例如:

首先将原始图像划分为值域块和定义域块两种大 小的块:值域块是完全无重叠的 n²小块 A1、A2、A3····, A.2; 定义域块是允许有部分重叠的 mi 大块 B.、B.、B., ...,B_{m²},(其中 m>n)。

然后寻找合适的仿射变换 ω 和定义域块 岛, 使得 $R_i = \omega_i(B_k)$,其中三维仿射变换 $\omega_i R^i \rightarrow R^i$ 的简化形式 为:

$$\omega \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} \mathbf{a} & \mathbf{b} & \mathbf{0} \\ \mathbf{c} & \mathbf{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{u} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} + \begin{pmatrix} \mathbf{e} \\ \mathbf{f} \\ \mathbf{g} \end{pmatrix}$$
(15)

在实际应用中,用式16的等价组合变换来代替式15寸 具有实现的可行性:

$$\omega_j = G_j \sigma \tau_j \sigma \Psi_j \tag{16}$$

其中,G,为灰度处理算子,包括比例因子 u 和灰度补 偿因子 g; τ, 为 1 个对折和旋转变换(通常 ι= ε);Ψ, 是 (x,y)平面上的紧缩变换,即将大小为 m²的 D, 映射成 大小为ni的块。这里

$$K_t = G_t \text{ or } n\Psi_t(D_t) \tag{17}$$

最后,只要找到了仿射变换 a, 和定义域块 D,,并

搜索引擎数据库也存在对搜索引擎的更新的问 题,距今越久远的记录,对系统决策的影响应该越小。 由于搜索引擎数据库内存贮的是每个搜索引擎对每个 关键词的查询能力,都是0~1之间的数值,数值越大表 示查询能力越强。因此通过定期让所有数值衰减,达到 更新数据库的目的。

结论 采用分布式的结构设计购书系统的搜索引 擎,各个节点搜索引擎之间的竞争能使它们不断完善 自己,并形成各自的特色,合作又使它们可以互相补 充、互相依赖。提出的数据库小型化、本地化等设计原 则注重于提高系统的性能与服务质量,特别是充分利 用了用户的反馈信息,可以使系统根据用户的不同需 求及时调整自身的运作,以提供最符合用户需求的购 书信息,从而提高了系统的服务质量。

参考文献

- 汪晓岩,胡庆生,李斌,庄镇泉,面向 Internet 的个性化智 能搜索引擎、计算机研究与发展、1999、36(9):1039~1046
- 邹涛,王继成,朱华字,金翔字,张福炎 WWW 上的信息挖 捆技术及实现. 计算机研究与发展, 1999, 36(8): 1019~ 1023
- 张卫丰,徐宝文, Web 搜索引擎框架研究, 计算机研究与 发展、2000、37(3):376~378
- 王军玲,赵沁平,一种基于类比的启发式搜索方法 计算机 科学,1998,25(5),33~37
- 常桂然,张晓辉. Web 信息检索服务系统与搜索引擎. 计 算机科学,1998、25(5):86~90
- 文栋辉,李光亚,赵振西,支持开放性的超媒体引擎,计算 机科学,1998,25(5):107~109

分别将其参数存储保存,待所有的值域子块都编码之 后,就可以完成运动图像的整体分形编码。

结束语 本文在文[2]的基础上介绍了一种运动 图像安全传输与存储平台的实现技术,重点研究了实 现该平台的压缩编码算法。该平台采用现代高速网络 技术,安全传输运动图像;采用国际上最先进的计算机 存储技术,智能化地长期保存运动图像信息;采用国际 上最新的软、硬件关键技术开发整个系统。其突出的优 点是可靠性高、体积小、成本低,可智能化地、长期不间 断地保存运动图像信息,而且不会发生影像失真或变 质,所以归档和管理容易,适合在我国大范围推广使 用。

参考文献

- 1 Tekalp A. M. Digital Video processing. Prentice-Hall International 1997
- remanonal (1997) 相泉, 邓锐、江立、肖明, 一种运动图像安全传输与存储平台的研究 华中师范大学学报(自然科学版),2001,8 钟玉琢, 乔秉新, 祁卫, 运动图像及其伴音通用编码国际标准——MPEG2. 清华大学出版社,1997 2
- 3
- 钟玉琢,冼伟铨,沈洪,多媒体技术基础及应用,清华大学 出版社 2000
- 品版社,2000 長乐南,数据压缩,东南大学出版社,2000 杨福生著,小波变换的工程分析与应用,科学出版社,2000