

# 网络多媒体教学资源主题搜索研究

杨仁广, 孟祥增

(山东师范大学 传播学院, 山东 济南 250014)

[摘要] 网络多媒体教学资源搜索与利用是信息教育中不可忽视的工作。基于主题搜索技术在专业领域中的应用, 建立教育主题词集、提取网络多媒体表征信息、改进主题搜索算法, 在已有的主题搜索器的基础上设计并实现了一个网络多媒体教学资源主题搜索系统, 用于搜索 Web 中的视频、音频、动画等多媒体资源, 为有效利用多媒体网络教学资源提供了一个良好平台。实验结果显示, 该系统能有效提高多媒体教学资源的搜索效率。

[关键词] 主题搜索; 教学资源; 多媒体; 主题搜索策略

[中图分类号] G40-057 [文献标识码] A

## 一、引言

网络多媒体教学资源是指存在于 Internet 中的多媒体教学资源。随着网络与多媒体技术的发展, Web 中的多媒体教学资源日益丰富, 成为教育领域的重要组成部分。如何快速、准确地找到特定主题的多媒体教学资源, 使其在信息化教育中充分发挥作用, 是教育技术工作者亟待解决的问题。20 世纪 90 年代发展起来的基于内容的多媒体检索技术,<sup>[1]</sup> 根据媒体的内容、语义、视觉特征以及媒体对象间的时空关系, 对多媒体数据库进行检索, 在一定程度上实现了搜索特定主题的多媒体。但基于内容的多媒体检索技术, 需要先建立起一个多媒体数据库, 并在多媒体数据库中对存储的多媒体进行内容、语义、视觉特征的特征, 造成整项工作耗时长、效率低, 甚至需要专门的机构负责多媒体教学资源的搜索工作。目前各大搜索引擎 (google、百度) 也相继推出了多媒体搜索引擎, 主要是利用网页中的相关文本提取描述多媒体信息的关键词进行多媒体信息检索, 这种搜索引擎能够直接、快速地从 Web 中寻找多媒体资源。但由于多媒体信息的内容极其丰富, 有时难以用几个简单的关键词进行描述, 甚至提取的关键词和多媒体内容不符, 导致搜索到的多媒体资源往往和教学主题没有关系。

基于主题搜索技术<sup>[2]</sup>在专业领域中的应用, 我们

对主题搜索的搜索过程、搜索算法以及搜索主题的判断进行改进, 设计并实现了一个网络多媒体教学资源的主题搜索系统, 用于搜索 Internet 中特定主题的多媒体教学资源。首先根据基础教育课本中目录、章节, 建立起各个学科的多媒体教学资源主题词集; 然后自动分析种子网站结构、获取网页链接、提取多媒体内容表征信息, 结合主题词集确定多媒体资源主题; 最后根据获取的网页链接进行下一个网页的分析, 最终实现对整个种子网站所有页面的分析判断, 进而寻找到此网站中包含的多媒体教学资源。下面就基于网络多媒体教学资源主题搜索系统的基础教育主题词集的建立、网络多媒体信息的内容表征与提取、系统的工作原理以及系统的主题搜索算法作较为详细的介绍。

## 二、基础教育网络多媒体教学资源主题词集的建立

为了确定在 Web 中搜索的多媒体资源的主题, 我们从人教版中小学课本中提取了与多媒体资源可能有关的主题词, 按学科、学段分类, 建立了高中语文、数学、物理、化学、生物, 初中语文、数学、物理、化学、生物, 小学语文、数学、科学、社会、思想品德与生活等 15 个主题词集, 同时每个词集下面又分为 3 个子词集: 视频词集、音频词集、动画词集。词集具体记

录数如表 1。不同学段、学科的主题词集中的词语有交叉,如表 2。

在网络多媒体教学资源主题搜索系统中,我们只搜索视频、音频、动画类多媒体资源。根据 CNNIC 发布的《第 19 次中国互联网络发展状况统计报告》,截至 2006 年底,中国网站的网页数量为 44.7 亿,其中文本和图像仍然是网页最主要的内容形式,分别占据

70.2% 和 29.5% 的比例; 视频网页占网页总数的 0.3%, 音频网页的比例几乎可以忽略不计, 而按照多媒体格式分类: swf 格式的网页占网页总数的 1%, MP3 格式的网页占网页总数的 0.1%。由此我们可以看出, 目前 Internet 上存在的网页大多以文本和图像形式存在, 所以此系统在搜索多媒体教学资源的时候图像资源暂时没有考虑。

表 1 基础教育网络多媒体教学资源主题词集记录数

类 型	小 学				初 中					高 中				
	语 文	数 学	科 学	社 会	语 文	数 学	物 理	化 学	生 物	语 文	数 学	物 理	化 学	生 物
视 频	165	42	72	49	138	49	84	40	155	124	27	76	117	60
音 频	90	18	68	73	180	21	22	37	19	152	35	52	19	28
动 画	308	285	63	56	149	114	94	40	120	166	72	98	75	41

表 2 网络多媒体教学资源主题词集中词语交叉示例

绿色食品 类型	分类	小 学				初 中					高 中				
		语 文	数 学	科 学	社 会	语 文	数 学	物 理	化 学	生 物	语 文	数 学	物 理	化 学	生 物
视频		√		√	√			√	√	√		√	√	√	
音频					√		√		√					√	
动画				√	√	√		√	√	√		√		√	

### 三、网络多媒体信息的内容表征与提取

Web 中的多媒体教学资源分布在各类网站中,基本上按三种方式存在:第一种作为网页的组成成分嵌入在网页中,我们称之为嵌入式多媒体,如 Flash、在线音视频播放等;第二种通过网页的锚文本链接,可以自由下载,我们称之为超链接式多媒体;第三种存在于多媒体网络的数据库中,允许检索、浏览,但常常需要账号和密码。其中以第一种、第二种形式存在的多媒体,在网页的源文件中都有相应的链接地址和表征此多媒体类型的格式扩展名。依据多媒体的格式扩展名,网络多媒体主题搜索系统可以判断出此多媒体类型。第三种由于是动态生成的,没有此类信息的提供,目前本系统对于这样的网页暂不分析。

《现代远程教学资源建设技术规范》<sup>[3]</sup>对远程教育中使用的多媒体教学资源属性从两个方面进行了描述:文件属性和内容特征。文件属性指文件名、文件格式、文件大小、文件建立时间、URL 等与文件相关的元信息;内容特征指多媒体本身包含的视听觉特征和结构、主题信息,如名称、构成、布局、画面颜色、对象形状、声强、音调等。其中内容特征的描述主要是人为主观描述,对于 Web 中海量分散的多媒体教学资源进行人工内容描述是不可能的。因此,考虑到 Web 中多媒体教学资源的特点和计算机智能信息处理能力,我

们基于“信息(认识论层次)是对事物属性表征”的思想,<sup>[4]</sup>提取包含多媒体的网页中对多媒体内容的描述信息。Web 中的多媒体通常以文件形式存在于网站服务器中,通过 HTTP 传输和网页一起呈现在客户端,网页中能够表征多媒体内容的描述信息有:网页 URL、网页 title、多媒体的文件名、超级链接文本、文件格式、URL 等。

网页一般是由超文本标记语言 HTML<sup>[5]</sup>(Hypertext Markup Language)编写的,其所包含的多媒体信息可以通过分析此网页的 HTML 标记获得。在网络多媒体教学资源主题搜索系统中,我们提取以下信息用来表征多媒体的主题:①网页的 URL;②网页<title>标签的文本内容;③网页中多媒体链接的 Anchor(锚文本);④多媒体链接的 URL。通过对大量包含多媒体的网页分析我们还发现:包含多媒体的网页链接在其父网页中通常以链接列表的形式出现,这些链接列表我们称之为“主题团”,<sup>[6]</sup>将“主题团”中包含的锚文本称之为“主题团”标题,它们对这些链接的主题起着指示性作用。为了对“主题团”的标题进行提前,提出 4 个启发性规则,每个“主题团”被限制在一对 table 标签内,并对内部嵌套的 table 标签进行合并。规则如下:

- (1) 该文本的字号比周围文本的大;
- (2) 该文本与周围文本的颜色不同;
- (3) 该文本字数很少(一般少于 10 个);

(4) 该文本独立成段。

如果满足其中任意两个,则将其认定为“主题团”标题。以此作为表征多媒体主题的重要信息。

#### 四、网络多媒体教学资源主题搜索系统的工作原理

网络多媒体教学资源主题搜索系统是专门为搜索 Web 中存在的多媒体教学资源而设计的一个主题搜索系统。主题搜索器的构造和常规的主题搜索器相同。主题蜘蛛是整个主题搜索器的核心,它首先从种子网站出发,通过 HTTP 协议请求下载种子网站页面,分析页面并提取链接,根据一定的启发式策略确定下一步的爬行方向,最终遍历种子网站中所有页面。主题搜索器从 Web 网页的文本信息中,高效准确地提取出多媒体教学资源内容表征信息,其检索的对象是文本信息,而不是多媒体内容本身。这里我们用文本作为多媒体内容的一种代理,其对多媒体的主题起着指示作用。图 1 是该主题搜索器的体系结构图,各个组成部分相互交错、相互依赖。

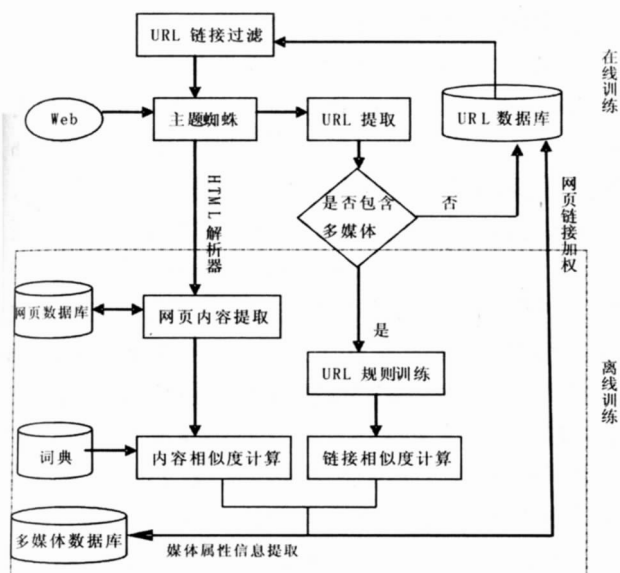


图 1 网络多媒体教学资源主题搜索器体系结构

整个搜索器分为在线训练和离线训练两部分(图中用虚线框起),其中在线训练主要负责网页信息提取和待爬行队列中 URL 的选择,离线训练主要负责 URL 规则训练和网页相关度的计算。其处理流程按照如下描述:

“主题蜘蛛”从互联网上抓取网页,提取网页的两部分信息:一是获取网页内容,以此来确定此网页与查询多媒体主题的相关度;二是提取网页链接,确定主题蜘蛛的即将爬行页面。“主题蜘蛛”通过 HTML 解析器获取此网页的文本信息,并将其和网页链接存入

“网页数据库”中。“网页内容相似度计算”用于判断此网页和查询多媒体主题的相关度。“主题蜘蛛”同时提取出网页的 URL,然后判断此网页是否包含多媒体,如果包含则进行“URL 规则训练”,将其用于“链接相似度计算”,经过内容和链接相似度计算后的网页链接与“URL 数据库”中的网页链接进行“网页链接加权”,从而确定下一步要爬行的网页,经过“URL 链接过滤”确定“网络蜘蛛”的爬行方向;如果不包含多媒体,则直接将提取的网页链接存入“URL 数据库”中,等待“网页链接加权”对其进行权值的分配。最终搜索得到的多媒体链接经过“媒体属性信息提取”后和表征多媒体内容的信息一起存入“多媒体数据库”中,同时也要将其存入“URL 数据库”中用于指示“主题蜘蛛”下一步爬行的方向。

#### 五、网络多媒体教学资源主题搜索算法

主题搜索算法是指专业搜索引擎为了提高在本专业领域中的搜索效率而提出的一个搜索算法,它能够即时、快速地为用户寻找到特定主题的内容。目前,该领域的很多专家、学者从理论和实践上做了很多研究工作,也提出了许多主题搜索算法,包括以 Shark-Search<sup>[7]</sup>和 Best-Fish<sup>[8]</sup>为代表的基于内容评价的搜索算法、以 PageRank<sup>[9]</sup>和 HITS<sup>[2,10]</sup>为代表的基于链接结构评价的搜索算法等,在实践中也取得了不错的效果,同时相应的专业搜索引擎也相继出现。但是应用于教育的专业搜索引擎现在还很少出现,专门对网络多媒体教学资源进行检索的搜索引擎更是没有。

我们基于网络多媒体教学资源在 Web 中分布的特点,对传统的 PageRank 和 Shark-Search 两种不同类型的主题搜索算法进行改进,将改进后的算法运用于网络多媒体教学资源的搜索中,实验结果表明搜索效率有显著提高。对 PageRank 和 Shark-Search 主题搜索算法的改进主要体现在两个方面:第一,内容相似度的计算方法;第二,链接相似度的计算方法,下边将作详细介绍。

##### 1. 内容相似度的计算

在前面我们提到“主题团”标题是描述多媒体主题的重要信息,在计算多媒体内容相似度的时候,我们把这个因素加入到计算过程中,具体如公式(1):

$$\text{Content\_score}(u_i) = \text{Score}(\text{block\_title}) [\beta * \text{Score}(\text{anchor}) + (1 - \beta) * \text{Score}(\text{url})] \quad (1)$$

其中,  $\text{score}(\text{block\_title})$  是链接  $u_i$  所在“主题团”



标题与主题的相关度,计算时采用向量空间模型 VSM。在向量空间模型中,所有的检索关键词  $t$  形成关键词集合  $T=(t_1, t_2, \dots, t_n)$ 。“主题团”标题文档  $D$  中的每一个文档  $d$  都被表示成一个范式的矢量  $V_i(d)=(t_1, w_1(d), \dots, t_n, w_n(d))$ ,其中  $w_i(d)$  为  $t_i$  在文档  $d$  中的权重,权重的计算采用 TF-IDF 词频统计方法计算,如公式(2)。采用向量空间模型计算“主题团”标题与主题的相关度,如公式(3)。Score(anchor)和 Score(url)分别表示链接  $u_i$  的锚文本和 URL 地址与主题的相关度,采用布尔模型进行计算; $\beta$  为相关因子,用以调节链接的锚文本和 URL 地址所占的比重。

$$w_i(d) = \frac{tf_i \log\left(\frac{N}{nt_i} + 0.01\right)}{\sqrt{\sum_{i=1}^n (tf_i \log\left(\frac{N}{nt_i} + 0.01\right))}} \quad (2)$$

其中:  $tf_i$  表示关键词  $t_i$  在文档  $d$  中出现的频率;  $N$  表示用于特征提取的全部训练文本的文档总数;  $nt_i$  表示出现关键词  $t_i$  的文档频率。

$$\text{Score}(\text{block\_title}) = \text{sim}(d, q) = \cos(\theta)$$

$$= \frac{\sum_{i=1}^n (w_i(d) * w_i(q))}{\sqrt{\sum_{j=1}^n w_j^2(d) * \sum_{j=1}^n w_j^2(q)}} \quad (3)$$

其中,  $w_i(q)$  为关键词  $t_i$  在查询  $q$  中的权重,通常当查询中包含就为 1; 否则就为 0。“主题团”标题与查询主题的相关度就表示为两个范化矢量之间夹角的余弦。

## 2. 改进 PageRank 算法

PageRank 算法是一种随机漫游模型,<sup>[6]</sup> 决定一个网页重要性的主要因素是指向该网页的链接个数。PageRank 算法在迭代计算过程中,权值是按当前网页的出度(out-degree)平均分配,没有考虑到网页的相对重要性。但在用户的实际访问过程中,用户会根据链接与主题的相似度,选择性地访问网页。例如一个页面  $X$  有 4 个链接,分别指向页面 A、B、C、D,4 个页面与主题的相似度分别为 0.2、0.4、0.6、0.8,则用户在选择网页连接时,选择页面 D 的机率要远远大于页面 A。

我们在利用 PageRank 算法对“待爬行队列中”的页面进行排序时,把网页的链接信息相似度加入到 PageRank 计算公式中。该算法认为:用户在查资料时,点击网页内链接的机率不是相等的,而是和两个因素——链接锚文本序列“主题团”的内容相似度和链接所指向网页的实际主题相关度(搜索过程中计算产生)有关。链接被点击的概率,同这两个因素成正

比。这里给出改进后的 PageRank 算法,如公式(4):

$$\text{Structure\_score}(u_i) = \frac{1-d}{N} + d * \sum_{i=1}^n \text{PR}(T_i) * P(T_i, u_i) \quad (4)$$

其中:  $\text{PR}(p)$  代表网页  $u_i$  的 PageRank 值;  $\text{PR}(T_i)$  代表网页  $T_i$  的 PageRank 值; 网页  $T_i$  指向网页  $u_i$ ;  $d$  为阻尼系数;  $P(T_i, u_i)$  为从页面  $T_i$  到达页面  $u_i$  的概率,计算方法如公式(5);  $N$  为已下载到待爬行队列中与主题相关的网页数量;  $n$  为链接到网页  $u_i$  的网页数量。

$$P(T_i, u_i) = \lambda * \frac{\text{score}(\text{block\_title})(u_i)}{\sum_{i=1}^n \text{score}(\text{block\_title})(i)} + (1 - \lambda) *$$

$$\frac{\text{sim}_{\text{link}}(u_i)}{\sum_{i=1}^n \text{sim}_{\text{link}}(i)} \quad (5)$$

其中:  $\sum_{i=1}^n \text{sim}_{\text{content}}$  表示从网页  $T_i$  中链出的所有网页内容相似度的集合;  $\sum_{i=1}^n \text{sim}_{\text{link}}(i)$  表示从网页  $T_i$  中链出的所有网页的实际主题相似度的集合;  $\lambda$  是一个影响因子,取值范围为 0~1。

## 3. 改进 Shark-Search 算法

Fish-Search 算法是 De Bra 等早期提出的一个主题网页动态爬行算法。Fish 算法将主题蜘蛛在 Web 中爬取网页的过程模拟为鱼群在大海中觅食的过程,当鱼找到食物时,鱼的繁殖能力就增强,反之鱼则逐渐消亡。Shark-Search 算法在 Fish-Search 算法的基础上做了两种主要的改进。首先,用一个连续的值函数来表示相关性,取值在 0~1 之间,而不是 Fish-Search 的二值判断;另外,待爬行链接的主题相关性受锚文本、锚文本上下文和父链接相关性继承的影响。文献<sup>[11]</sup>对网页中不同区块的链接进行聚类,然后将相同类的所有链接锚文本作为该类的描述文本,用来替代 Shark-Search 算法中锚文本上下文对链接相关性的影响。文献<sup>[12]</sup>从网页页面、链接块以及链接本身 3 个粒度上对网页的相似性分别进行计算,然后将三者按照不同权重结合,进而确定整个网页的相似性。

在链接结构方面,我们发现包含多媒体的主题网页表现出“资源相邻性”的特点。所谓“资源相邻性”是指在一个网站中,包含的多媒体资源往往存在于这个网站的某一部分或某几部分区域中,并且处于同一区域的多媒体资源的主题也是相同的。根据此特点我们做出如下假设:

(1) 如果一个网页是与主题相关的包含多媒体的网页,那么此网页的子链接很可能是与主题相关的包

含多媒体的网页;

(2)如果一个网页是与主题相关的包含多媒体的网页,那么此网页在父网页中的兄弟链接很可能是与主题相关的包含多媒体的网页。

所以在计算网页链接相关度时,我们用父网页和兄弟网页的链接相关度来揭示链接结构对一个 URL 链接相关度的影响。为了把这种影响实时地反馈给每个子链接,引入一个动态因子。用公式(6)来表示链接结构对一个 URL 链接相关度的贡献:

$$\text{Structure\_score}(\text{block\_title})(u_i) = \sum_{j \in \text{兄弟}} \lambda(d_j) P(d_j) / t \quad (6)$$

其中: $u_i$ 是正在爬行的链接, $t$ 是父链接的总数, $\lambda(d_j)$ 是动态因子,用公式(7)进行计算; $P(d_j)$ 表示从父链接继承来的链接相关度和已爬行过兄弟链接的平均链接相关度,用它来衡量通过父链接能爬行到多少主题相关页面的能力,用公式(8)进行计算。

$$\lambda(d_j) = (n' + \theta) / (n + \theta) \quad (7)$$

其中: $n'$ 是父链接 $d_j$ 的已爬行子链接中主题相关页面的个数; $n$ 表示父链接 $d_j$ 已爬行子链接的总个数; $\theta$ 是归一化因子,通常取0.5。在爬行的过程中, $\lambda(d_j)$ 会不断地动态调整。

$$P(d_j) = (1 - \sigma) R(d_j) + \sigma \sum_{k \in \text{兄弟}} P(d_k) / N \quad (8)$$

其中: $\sigma$ 是偏置因子, $R(d_j)$ 为父链接 $d_j$ 的主题相关度, $d_k$ 为 $d_j$ 的一个已爬行子链接, $N$ 为 $d_j$ 已爬行子链接的总数, $\sum_{k \in \text{兄弟}} P(d_k) / N$ 是父链接 $d_j$ 中已爬行子链接的平均链接得分。

#### 4. 内容相似度和链接相似度的归一化

为了提高整个网页的主题相关性和权威性,我们采用内容相似度和链接相似度按不同权值相加所得

结果来标志。在这里将二者归一化,计算得到的值作为“主题蜘蛛”即将爬行链接的依据。计算公式为:

$$S_i = \sigma * \text{Content\_score}(u_i) + (1 - \sigma) * \text{Structure\_score}(u_i) \quad (9)$$

## 六、实验结果

### 1. 参数的选择和评价指标

经反复测试,对于改进的 PageRank 算法,选择  $d=0.85, \lambda=0.5, \sigma=0.6$  时,实验效果最好;对于改进的 Shark-Search 算法,选择  $\beta=0.85, \sigma=0.6, \lambda=0.5$  时,实验效果最好。用查准率和查全率两个指标来评价算法的效果。

### 2. 仿真试验环境

两种改进的算法用 Java 语言实现。为提高搜索效率,采用多线程技术同时搜索不同的站点,系统开启了 10 个线程。实验环境为:Windows XP 操作系统 PIII CPU, 512M 内存。

### 3. 实验结果

人工选择 10 个物理教育网站作为种子(6 个包含嵌入式多媒体,4 个包含超链接式多媒体),为了更好地检验搜索结果,我们将初中物理和高中物理主题词集合并为统一的物理词集(192 个词条)先用通用搜索算法运行,结果如表 3 所示;然后分别用标准 Fish 算法、改进的 PageRank 算法和改进的 Shark-Search 算法运行,结果如表 4 所示。

表 3 通用搜索算法实验结果

种子个数	网页总数	有效网页个数	有效网页占有率	运行时间	平均爬行速度
10	73952	3059	4.464%	30.91	40 个/分钟

实验还从算法的查准率和爬行时间的关系出发,对三种算法进行测试,每隔 1 小时统计一次,结果如表 5 所示。

表 4

三种算法实验结果比较

算法	种子个数	网页总数	有效网页个数	运行时间	平均爬行速度	查准率	查全率
标准 Fish 算法	10	36872	984	18.44 小时	33 个/分钟	2.67%	32.17%
改进的 PageRank 算法	10	21691	1643	24.51 小时	14 个/分钟	7.58%	53.71%
改进 Shark-Search 算法	10	15892	2048	15.27 小时	17 个/分钟	12.89%	66.95%

表 5 三种算法查准率随时间变化实验结果比较

查准率/(%)			时间/小时
标准 Fish 算法	Shark-Search 算法	改进 Shark-Search	
2.256%	7.485%	6.527%	1
3.246%	12.267%	16.463%	2
3.568%	19.619%	21.493%	3
3.896%	21.572%	25.691%	4

### 4. 结果分析

由表 4 可以看出,改进的 PageRank 算法和改进的 Shark-Search 算法由于需要计算内容相似度和链接相似度,在平均速度上低于标准的 Fish 算法。改进的 PageRank 算法在运行的过程中,计算网页链接相似度时要计算网页中超链接的爬行概率在平均速度上要低于改进的 Shark-Search 算法。两种改进的算

法在查准率和查全率方面比通用搜索算法都有较大的提高。虽然表 4、表 5 显示改进的 PageRank 算法比改进的 Shark-Search 算法在搜索精度和搜索效率上都有一定差距,但主要原因是两种算法的原理不一样,适用的网页也不同。

### 七、结束语

以上详细介绍了整个网络多媒体教学资源主题搜

索系统的各个环节,同时我们也将系统搜索的结果用于本实验室开发的 Web 多媒体资源搜索系统(<http://www.cbxy.sdnsu.edu.cn/cbxy/WebRetrieval/index.asp>)中,实验效果良好。今后我们将做以下工作:①扩展基础教育主题词集,扩大搜索的范围,即时更新“多媒体数据库”中的记录;②继续提高多媒体主题搜索算法的效率,着重优化实验算法中各个参数;③考虑对存在于多媒体网络的数据库(动态网页)中多媒体资源的获取。

### [参考文献]

- [1] 章毓晋[J].基于内容的视觉信息检索[M].北京:科学技术出版社,2003.
- [2] Bharat K,Henznger M R.Improved Algorithms for Topic Distillation in a Hyperlinked Environment[C]//Proceedings of SIGIR Conference on Research and Development in Information Retrieval New York,1998:104~111.
- [3] 教育部信息化技术标准委员会.现代远程教学资源建设规范[S].2001.
- [4] 孟祥增.多媒体网络教学资源的内容特征提取与搜索研究[J].电化教育研究,2007,(12):33~37.
- [5] 杨仁广,孟祥增.一种基于网页内容和链接分析的主题搜索算法[J].情报杂志,2008,(6):64~66.
- [6] 宋宇,孟祥增.基于改进 Fish-Search 算法的多媒体检索[J].计算机工程,2008,34(11):189~193.
- [7] Bra D P,Houben G, Kornatzky et al. Information Retrieval in Distributed Hypertexts[C].Proceeding of the 4th RIAO Conference,1994.
- [8] Cho J,Garcia-MolinaH,Page L.Efficient Crawling Through URL Ordering[J].Computer Networks,1998,30(1~7):161~172.
- [9] 曾春,邢春晓.基于内容过滤的个性化搜索算法[J].软件学报,2003,14(5):999~1004.
- [10] Menczer F. Complementing Search Engines with Online Web Mining Agents[J].Decision Support Systems,2003,35(2):195~212.
- [11] 苏祺,项锴,孙斌.基于链接聚类的 Shark-Search 算法[J].山东大学学报(理学版),2006,41(3):1~4.
- [12] 陈骏,陈竹敏.基于网页分块的 Shark-Search 算法[J].山东大学学报(理学版),2007,42(9):62~66.

## DRA 荣升国标 大洋力助其产品化

4月19日,《多声道数字音频编解码技术规范》(DRA)国家标准发布会在钓鱼台国宾馆隆重举行,发布会由国家质量监督检验检疫总局、工业和信息化部与广东省人民政府联合召开。作为参与其产品化、产业化进程的唯一广电厂商,中科大洋公司应邀参加此次发布会,并展出了应用 DRA 标准的大洋 D3-Edit 系列非编产品,受到了与会者的广泛关注。

DRA 是国内首个,同时也是除杜比、DTS 之外的世界第三个环绕声伴音标准,拥有自主知识产权;该标准继 2007 年被原信息产业部颁发为电子行业标准后,又于今年 2 月被批准为数字音频编解码技术国家标准。在此次发布会上,来自国家质检总局、工业和信息化部、广东省政府的多位领导亲临现场,国家标准化委员会副主任方向先生宣布《多声道数字音频编解码技术规范》(DRA)国家标准正式发布!

在 DRA 标准的应用过程中,大洋公司充分彰显了勇于接受新事物和支持民族音频技术发展的精神,力助其推广应用;DRA 标准诞生以后,公司迅速组织技术人员开展应用研究,经过一段时间的努力,目前,大洋公司的 D3-Edit 系列非编产品、磐石视频服务器等产品已全面支持 DRA 标准;除产品应用外,大洋公司还对 DRA 标准的修订给予了有益的建议,并通过积极参加国家标准汇报会等形式,与业界同仁共同推动 DRA 标准的产业化和推广应用。

DRA 上升为国家标准,意味着国外公司长期垄断中国音频技术市场的局面被打破,中国企业在音频领域有了属于自己的标准;而大洋公司在 D3-Edit 等产品中积极应用 DRA 标准,则标志着其在数字电视节目制作环绕声伴音技术上的领先地位进一步确立!