

# JZSearch大数据精准搜索关键技术

JZSearch Big Data Precise Search Key Technology



张华平 博士 副教授 大数据搜索与挖掘实验室 kevinzhang@bit.edu.cn @ICTCLAS张华平博士 2016.10









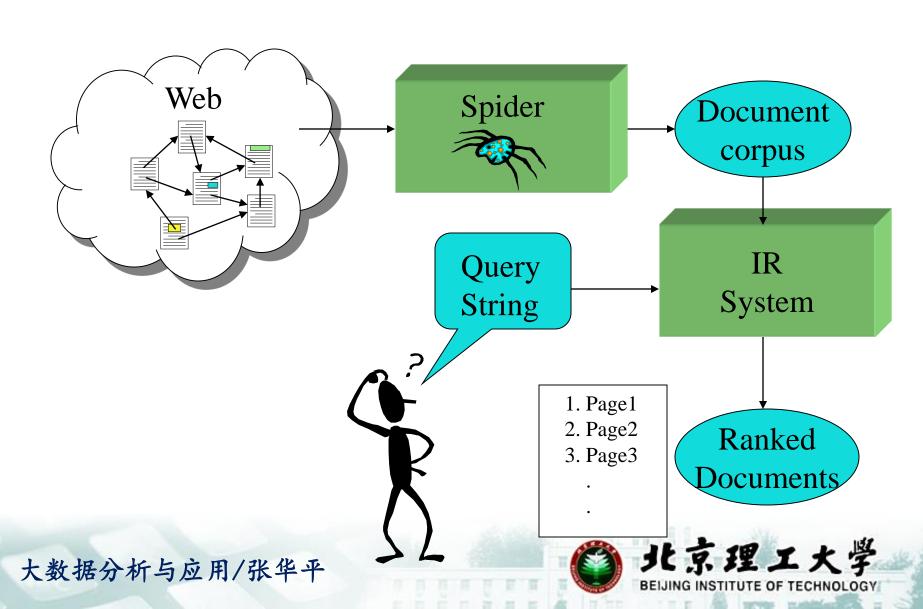


I 简单Web搜索引擎基本原理

- 当前搜索技术的主要不足
- IJZSearch大数据精准搜索引擎
- JZSearch大数据搜索应用案例



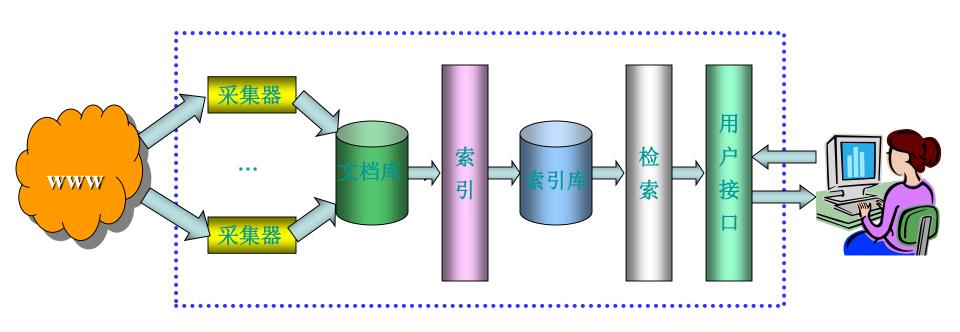
#### Web Search Engine Using IR



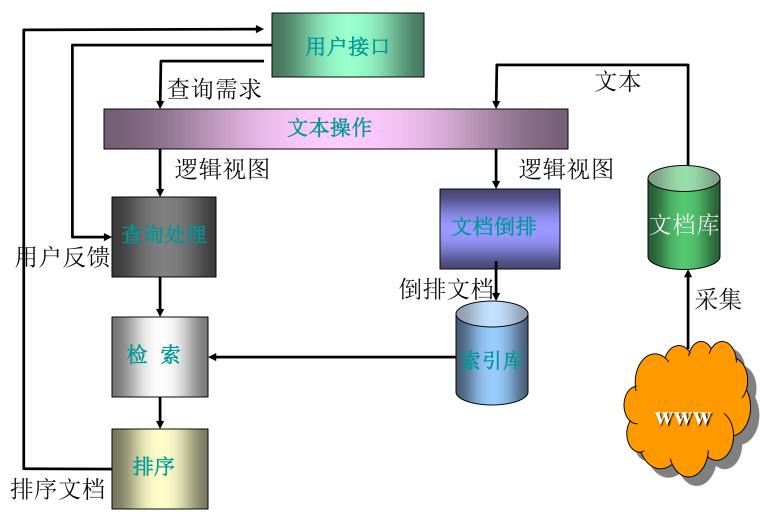


# 最简单的搜索引擎

#### 7搜索引擎结构





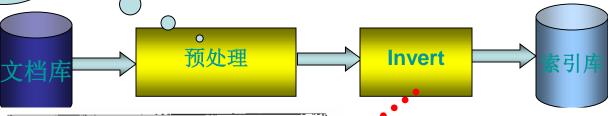


大数据分析与应用/张华平



# 倒排索引

文档分析,编码 识别,词语切分, 去停用词等,



Document	Text
1	Pease porridge hot, pease porridge cold,
2	Pease porridge in the pot,
3	Nine days old.
4	Some like it hot, some like it cold,
5	Some like it in the pot,
6	Nine days old.

文档倒排,生成Inverted Files

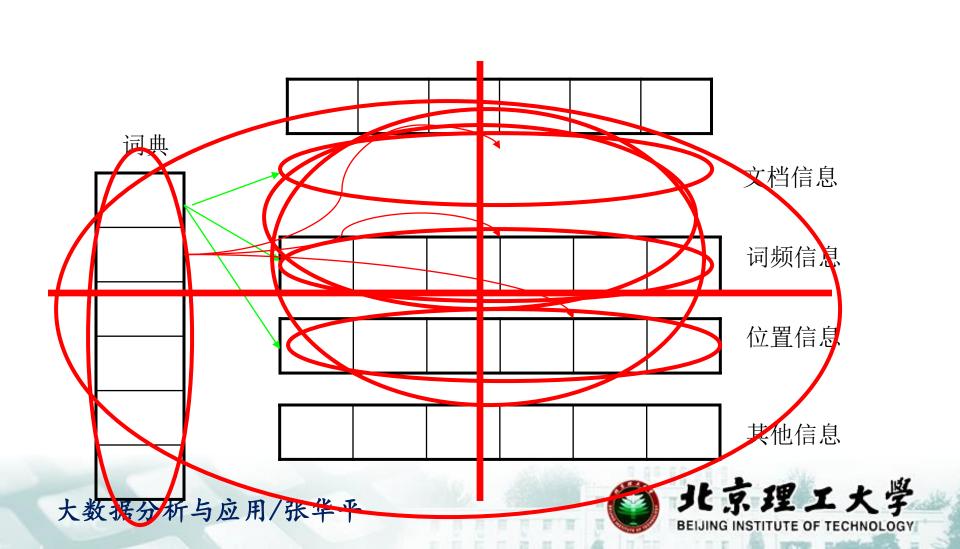
Number	Term	(Document; Words)
1	cold	(1; 6), (4; 8)
2	days	(3; 2), (6; 2)
3	hot	(1; 3), (4; 4)
4	in	(2; 3), (5; 4)
5	it	(4; 3, 7), (5; 3)
6	like	(4; 2, 6), (5; 2)
7	nine	(3; 1), (6; 1)
8	old	(3; 3), (6; 3)
9	pease	(1; 1, 4), (2; 1)
10	porridge	(1; 2, 5), (2; 2)
11	pot	(2; 5), (5; 6)
12	some	(4; 1, 5), (5; 1)
13	the	(2; 4), (5; 5)

に学



#### 倒排索引结构

7 索引文件结构:不管怎么变化基本都由这几部分组成





# 倒排索引的挑战

对Fast:如何快速构造倒排索引?

7Fast:如何构造索引使检索尽可能的快?

Minimized: 如何使索引尽可能小?

70nline Indexing: 如何构造动态文档集的索引(增量, 差量索引和索引更新)?

→ Scalable:如何在资源有限的情况下构造海量数据的索引?





信息检索

集合论模型

模糊集合论模型扩展布尔模型

代数模型

广义向量模型 潜语义标引模型 神经网络模型

概率模型

推理网络模型 信任度网络模型

经典模型

布尔模型 向量模型 概率模型

结构化模型

非重叠链表模型 邻近节点模型

浏览

扁平示模型 结构导向模型 超文本模型

大数据分析与应用/张华平

检索:

过滤

用户任务

特别检索

浏览

信息检索模型分类



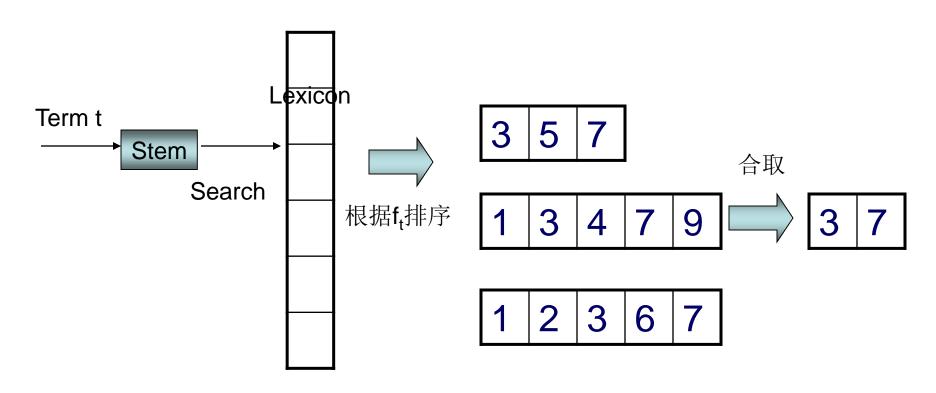
北京理工大学

#### 布尔查询

- 7一种简单的检索模型,建立在经典的集合 论和布尔代数的基础上。
- 7 遵循两条基本规则: 每个索引词在一篇文档中只有两种状态: 出现或不出现, 对应权值为0或1。
- 力查询是由三种布尔逻辑运算符 and, or, not 连接索引词组成的布尔表达式。



### 合取布尔查询处理





# Ranking和信息检索

- 2 经典布尔模型能精确判断文档是否出现某一查询,但并不能给出相关性排序
- ↑ 信息检索是一个查询Q和文档Dd相似度计算过程:

$$M(Q,D_d) = Q \bullet D_d = \sum_{t=1}^n w_{q,t} \bullet w_{d,t} = \sum_{t \in O} w_{q,t} \bullet w_{d,t}$$

7 存在一个问题: 当Q包含常用词t时, 那些包含比较多t的文档总是排在前面, 其他的非常用词根本不起作用, 所以需要根据inverse document frequency (IDF)计算Term的权重 w<sub>t</sub>:

$$w_t = \frac{1}{f_t}$$



 $M_{2}$ 



 $\vec{x}$   $\overline{M_1M_2}$ 

- 7 向量(矢量, vector): 既有大小又有方向的量, 通常用有向线段表示, 记作 或者
- 7考虑从空间坐标系原点出发(其他向量可以· $M_1$  平移到原点出发)的向量  $\vec{x}$  ,终点坐标为  $\langle x_1, x_2, \cdots, x_n \rangle$  ,我们称之为一个n维向量
- 力向量的运算: 加、减、倍数、内积  $\vec{x} \pm \vec{y} = \langle x_1 \pm y_1, x_2 \pm y_2, ..., x_n \pm y_n \rangle$   $\lambda \vec{x} = \langle \lambda x_1, \lambda x_2, ..., \lambda x_n \rangle$





#### 向量的模、距离和夹角

$$|\vec{x}| = ||\vec{x}|| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

7向量的模(大小):

7向量的(欧氏)距离

$$dist(\vec{x}, \vec{y}) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

**7**夹角α

$$\cos \alpha = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \times |\vec{y}|}$$





#### 向量空间模型

- 向量空间模型(Vector Space Model, VSM)是康奈尔大学 Salton等人上世纪70年代提出并倡导. 原型系统SMART\*
- 7 term独立性假设: term在文档中的出现是独立、互不影响的。
- 力查询和文档都可转化成term及其权重组成的向量表示,都可以看成空间中的点。向量之间通过距离计算得到查询和每个文档的相似度。



#### 文档-标引项矩阵(Doc-Term Matrix)

n篇文档,m个标引项构成的矩阵 $A_{m*n}$ ,每列可以看成每篇文档的向量表示,同时,每行也可以可以看成标引项的向量表示。

$$A_{m*n} = \begin{bmatrix} d_1 & d_2 & \dots & d_n \\ t_1 & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \dots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$







- **⊅** 查询q: (<2006,1>,<世界杯,2>)
- **7** 文档*d*<sub>1</sub>: (<2006,1>,<世界杯,3>,<德国,1>,<举行,1>)
- 7 文档d<sub>2</sub>: (<2002,1>,<世界杯,2>,<韩国,1>,<日本,1>,<举行,1>)

	d	d	2	q
2002	0	1		$\lceil 0 \rceil$
2006	1	0		1
世界杯	3	2		2
德国	1	0		0
韩国	0	1		0
日本	0	1		0
举行	_1	1_		$\lfloor 0 \rfloor$



# 一个例子(续)

#### 7查询和文档进行向量的相似度计算:

■采用内积:

**①**文档 $d_1$ 与q的内积: 1\*1+3\*2=7

**①**文档 $d_2$ 与q的内积: 2\*2=4

■夹角余弦:

●文档d2与q的夹角余弦:

$$\frac{7}{\sqrt{12\times5}}\approx0.90$$

$$\frac{4}{\sqrt{5\times8}} \approx 0.63$$



北京理工大学



# 典型的静态信息检索

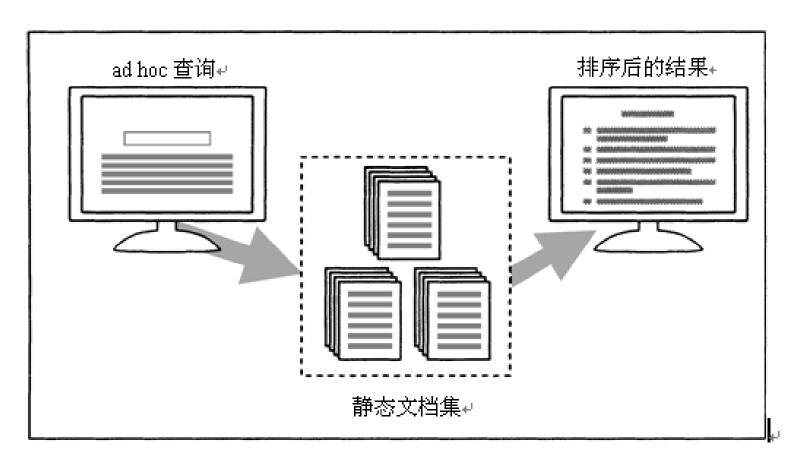
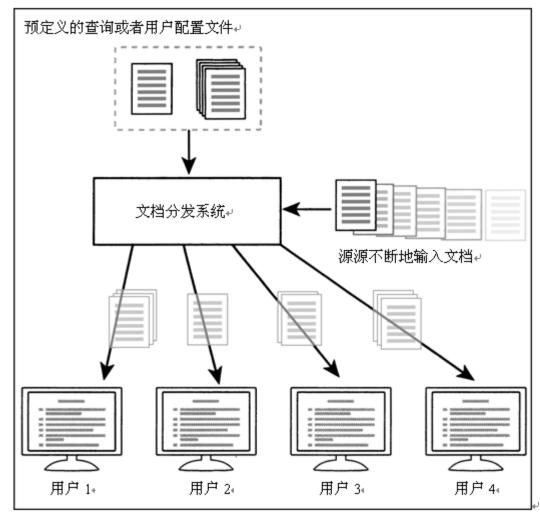


图 1-1 文档检索→





# 典型的动态信息流过滤过程





### 信息检索vs.信息流过滤的差异

#### 7文档集:

■规模有限,相对静态 vs. 规模巨大、动态变化

#### 7用户需求:

■多变多样,海量用户 vs. 相对稳定、用户受限

#### 7技术方法:

■对文档进行倒排索引并检索 vs. 查询生成自动机,进行快速扫描过滤



# 信息检索vs.信息流过滤的共同

7同一理论支持;

7相同的支持技术:自然语言处理、机器学 习、大规模架构设计

7小规模的信息安全可以直接规约为信息检索问题。





#### 信息检索需研究的五大问题

- 7搜索排序问题:信息检索模型,解决搜索的 根本问题;
- 7效果提升问题:信息检索策略,提高效果的 策略与技巧;
- ▶性能问题:倒排、索引压缩;
- 7架构问题:并行、分布式、云计算;
- 7 拓展问题:多语言、跨语言等。







简单Web搜索引擎基本原理

Ⅱ 当前搜索技术的主要不足

- IJZSearch大数据精准搜索引擎
- JZWearch大数据搜索应用案例



# 当前不足初探

#### 7信息重复冗余

普通搜索

关键字: 白蜡 搜索类型: 按全文搜索

1文示天空・ 技主义技系 \*

行业: -请选择-

栏目类别: 全部类别

地区: -地区- ▼

日期: 近一月

搜索

<b>当前位置:</b> 中国搜标网 >>普通搜索		共16条记录 没有	有查到所需的信息 <b>?</b>
信息标题	地区	日期	收藏
关于驻马店市园林管理局绿地养护项目采购公告	河南省	2012-09-17	豐收藏
关于驻马店市园林管理局绿地养护项目招标公告	河南省	2012-09-17	豐收藏
询价公告(绿化苗木、12-1 <b>4-37</b> )	新疆	2012-09-17	豐收藏
神新能源公司2012年秋季绿化工程招标公告	新疆	2012-09-14	豐收職
滦河东路绿化景观改造工程(六标段)滦河东路绿化景观改造工程(六标段)招	河北省	2012-09-11	豐收藏
滦河东路绿化景观改造工程(六标段)招标公告	河北	2012-09-11	[豐收藏]
黄骅港联轴节、维修车间维修工具采购-采购公告	河北	2012-09-07	豐收藏
兴辰道道路及配套管线工程监理招标公告	天津市	2012-09-07	豐收藏
宜白路道路及配套管线工程监理招标	天津市	2012-09-07	豐收藏
宁夏银川市贺兰县农业综合开发办秋季造林苗木竞争性谈判招标公告	宁夏	2012-09-06	豐收藏
宁夏银川市贺兰县农业综合开发办秋季造林苗木竞争性谈判招标公告	宁夏	2012-09-06	豐收藏

大数据分析与应用/张华平



北京理工大学

BEIJING INSTITUTE OF TECHNOLOGY



## 当前不足初探

#### 7白腊-石蜡 是专业同义词, 缺乏专业知识 关联



关键字: 石蜡 搜索类型: 按全文搜索 ▼ 栏目类别: 全部类别 ▼

地区: -地区- ▼ 行业: -请选择- ▼ 日期: 近一月

搜索

<b>当前位置:</b> 中国搜标网 >>普通搜索		共32条记录 <b>没有</b>	查到所需的信息?
信息标题	地区	日期	收藏
广州市番禺区政府采购医疗设备项目招标公告	广东	2012-09-21	豐收藏
天津市南开医院全自动组织脱水机等设备采购项目招标公告	天津市	2012-09-21	豐收藏
巴中市中心医院一批医疗设备采购招标公告	四川省	2012-09-20	□□收藏〕
四川省: 泸州医学院教学、科研设备第二批采购项目(第二次)中标公告	四川省	2012-09-20	□□收藏〕
北京市密云水库医院病理科医疗设备购置政府采购项目招标公告	北京市	2012-09-17	豐收藏
北京市密云水库医院病理科医疗设备购置政府采购项目招标公告	北京	2012-09-17	豐收藏
盐亭县电化教育馆教学仪器采购招标公告	四川省	2012-09-13	□□收藏]
四川省: 泸州医学院教学、科研设备第二批采购项目(第二次)中标公示	四川省	2012-09-12	□□收藏]
四川省人民医院试剂比选通知	四川省	2012-09-12	□□收藏〕
9月份超期老品	山东省	2012-09-11	豐收藏



# 当前不足初探

#### 7搜索禁用语禁而不止,搜索结果完全无关



市场热线: 010-82744908 82744968 82743201 **合作咨询:** 010-82744228
Copyright © 2005-2009 中国搜标网 All Rights Reserved. 京ICP证070104号 北京市公安局海淀分局备案编号1101081890



### 通用搜索 vs. 大数据专业搜索

- · 召回率不是通用搜索的考量范围,只关注P@10,P@30:
  - 体量大,一大遮百丑
  - 查询:中国证监会所有的负面信息;
- 体量大, 但返回给用户的有限;
  - 返回给用户的信息不足2000条
  - 查询: 新疆极端势力的动态情报;
- 专业垂直整合缺乏
  - 专业业务的逻辑不在通用搜索的业务范畴内;
  - 中关村东路80号5楼1203室的邮编是多少?联想集团的信息 汇总





# 关于作业: 最近三年(ACL, SIGIR,CIKM,WWW, SMP) 最经典的Tutorial



大数据 搜索

I 简单Web搜索引擎基本原理

当前搜索技术的主要不足

- IJZSearch大数据精准搜索引擎
- JZSearch大数据搜索应用案例



#### JZSearch大数据精准搜索

#### • 搜索基本功能:

• 多字段关联搜索、指定字段排序、精确搜索与模糊搜索

#### • 搜索特色功能:

- 内嵌正负面情感等极性分析、语义联想搜索、临近搜索、搜索结果去重;
- 内嵌了ICTCLAS智能分词系统:
- 数据库实时同步:数据库增删改10秒内即可同步到搜索;

#### • 搜索维护功能:

• 单点故障容错;支持增量索引;自动备份与恢复机制;自动缓存机制;自动优化机制;搜索屏蔽与恢复;



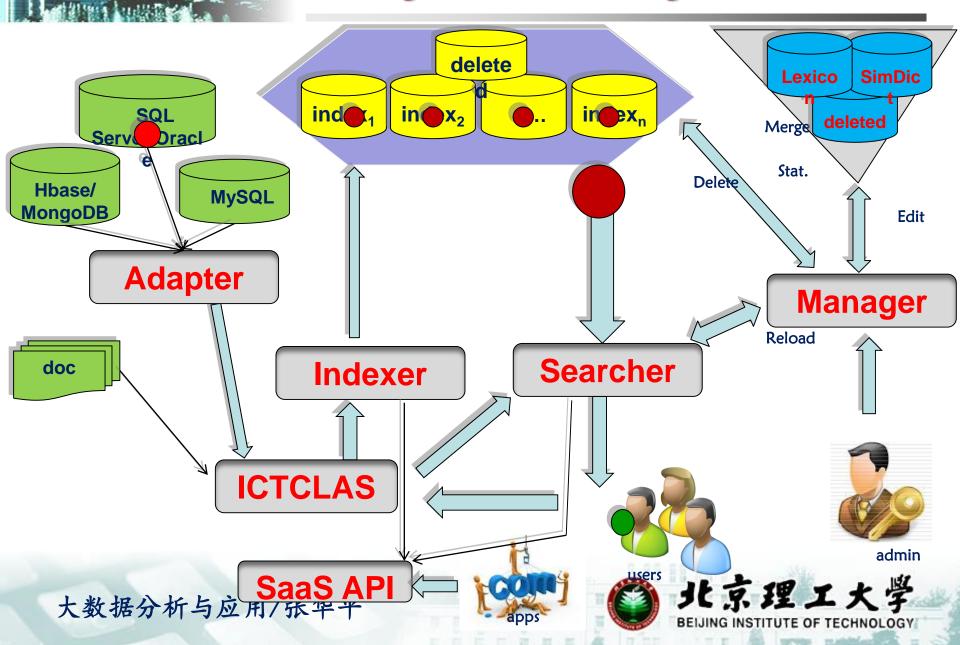


#### JZSearch内部语法示例

- **7** [FIELD] \* [NEAR] 尚福林 ##负面JZSearch## 12
- **↗** [FIELD] price [RANG] 1.0 9.0 [FIELD] name [AND] 牛奶儿童
- **7** [FIELD] name [PREF] 张
  - 姓名字段name必须以"张"作为前缀开头
- **↗** [FIELD] id [PREC] 123
  - 字段id必须以"123"精准匹配,如"1234"或者 "0123"均不作为匹配结果;
- **オ** [field] content [complex] 统计局||中国统计局||CPI 骗人|| 砖家 10

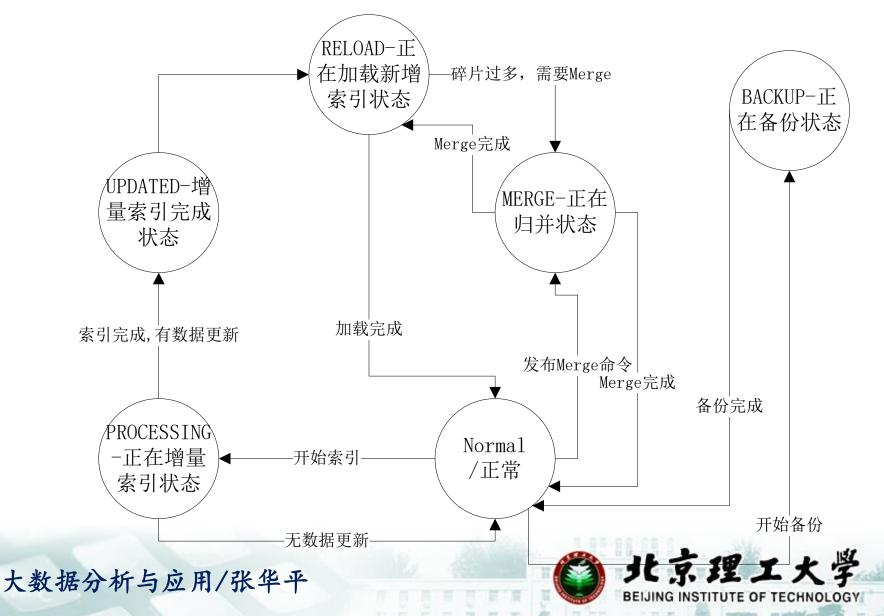


#### JZSearch Big Data Search Engine Architecture





## JZSearch 状态转移图





#### 索引压缩技术

- 7 倒排索引中倒排表的大小可以通过齐普夫分布来估计。
- 7 它描述了自然语言中每个词的词频分布情况:如果把词出现的频率按由大到小排列,那么每个词的词频与它排序序号的乘积是一个常量。
- 7 常量取1时齐普夫分布中前5个词的词频分布如下:

排序序号。	频率。	常里。
1.5	1.00.1	1.5
2.1	0.50.,	1.5
3.1	0.33.,	1.5
4.,	0. 25.,	1.,
5.,	0. 20.,	1.,



#### 索引压缩:哈夫曼编码

7 在英文中, e的出现概率很高, 而z的出现概率则最低。当利用哈夫曼编码对一篇英文进行压缩时, e极有可能用一个位(bit)来表示, 而z则可能花去25个位(不是26)。用普通的表示方法时, 每个英文字母均占用一个字节(byte), 即8个位。二者相比, e使用了一般编码的1/8的长度, z则使用了3倍多。倘若我们能实现对于英文中各个字母出现概率的较准确的估算, 就可以大幅度提高无损压缩的比例。



# O(Q(tf<sub>max</sub>))

# JZSearch索引压缩

### a 字节对齐压缩(Byte-Aligned)

- 0-63 00xxxxxx
- 64-(16K-1) 01xxxxxx xxxxxxx
- 16K-(4M-1) 10xxxxxx xxxxxxx xxxxxxx
- 4M-(1G-1) 11xxxxxx xxxxxxxx xxxxxxx xxxxxxx
- **0** 00000000
- **1** 00000001
- **...** ...
- **6**3 00111111
- **64** 01000000 01000000
- **65** 01000000 01000001
- 7 通过使用字节边界来实现BA压缩,运行时会付出些许代价,但可以获得一定的压缩比。BA算法易于实现,而且压缩比比较高(使用停用词后,压缩文件大小占未压缩索引文件大小的比例大约是

大数据分析与应用/张华平



北京理工大学





### 7 固定长度的索引压缩示例

- 7 给定任意一个索引词t1,我们考虑其索引的条目,假定t1在文档1、 文档3、文档7、文档70和文档250中出现。
- 7 BA压缩使用最高的两个比特位来表明存储该值需要的字节数。对于前4个值,只需要一个字节就可以表示;对于最后一个值——180,需要两个字节来表示。需要注意的是:我们只需要计算倒排表中每一项的差。最后的差值为:250-70=180。我们需要计算所有的最终数值。这些数值和压缩后的相应比特字符串见下表。(字节对齐方

式压缩)	值。	压缩的比特串。
	1.1	00 000001.,
	2.1	00 000010.,
	4.,	00 000100.,
	63.,	00 111111.
	180.,	01 000000 10110100.





北京理工大学



### 7 固定长度的索引压缩示例

在没有压缩之前,倒排表中的每条记录需要4字节,5条记录总共需要20个字节。这些数值和它们对应压缩前的比特字符串详见下表。 (基准线:压缩前)

值。	未压缩的比特串。
1.1	00000000 00000000 00000000 00000001.,
3.,	00000000 00000000 00000000 00000011.
7.1	00000000 00000000 00000000 00000111.
70.,	00000000 00000000 00000000 01000110.,
250.,	00000000 00000000 00000000 11111010.,

7 结论:在这个例子中,未压缩的数据需要160比特,使用BA压缩后仅需要48比特。

大数据分析与应用/张华平

### 索引词的处理

- 7 索引词典动态设计:
  - 优点:可以快速索引新词
  - 缺点: 动态词典组织效率远低于静态词典
- 7 JZSearch采用静态词典, 主要包含:
  - 汉语常用词
  - 英文常用词: 词形自动转换
  - 英文未知词表示: n-gram拆分
  - 数字组合: n-gram拆分





I 简单Web搜索引擎基本原理

Ⅲ 当前搜索技术的主要不足

- III JZSearch大数据精准搜索引擎
  - JZSearch大数据搜索应用案例

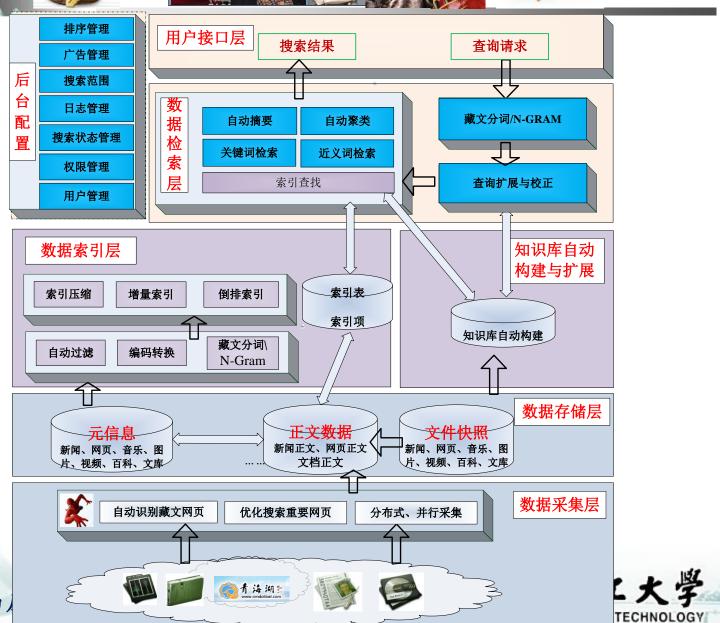




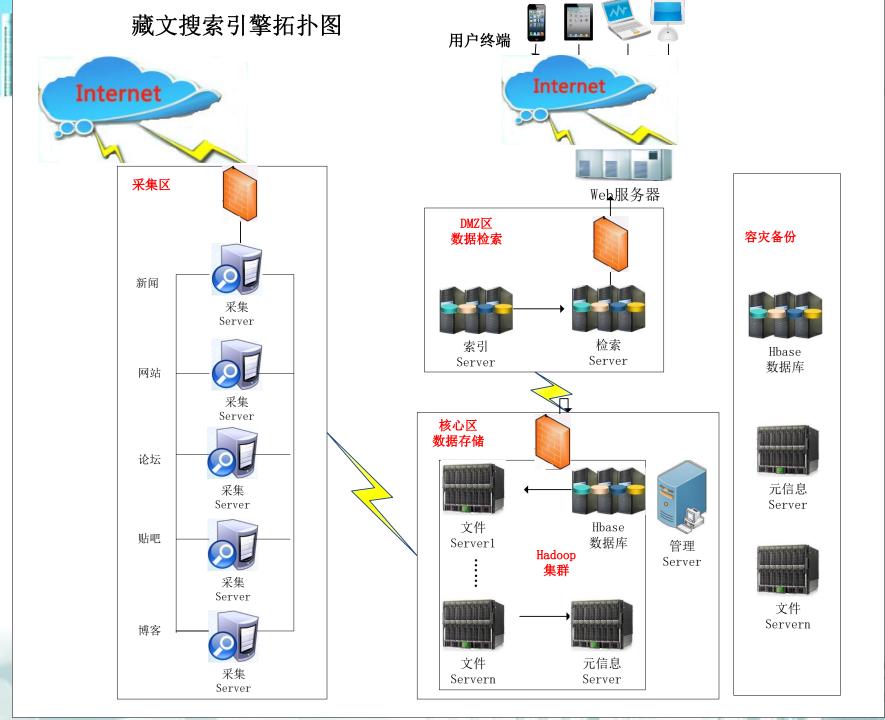




### 软件架构



大数据分析与







## 搜索引擎-爬虫技术



网络舆情挖掘系统

网络舆情挖掘系统

30完成

21 故宫建福宫袖指成富豪会所 工作人员:不能留宿

观海知心沅

新华社北京分社

2011-05-13 09:52

840M 🜒 🚵 🖳 🤄 🔵 🔞 0 🔁 缩放:100% 日报 杳看

来源:微博相册- 评论(0)- 转发(0)-发布时间: 2011-12-03 11:56:50



# 公司搜索

HAZIMAN DINEH	BOX W YO	1 11. 78 DE-F	49-21-412 A . 1							
ZsearchAgent	tClient								<u> </u>	
	127 . 0 . 0 \IndexFile\field		端口号: 8001		搜索编	号: 0 默认值		毎页结果数: 显示格式:	20	•
搜索结果:	美国恩艾仪器	8有限公司				搜索	· 下一页	上一页	过 退	ť
搜索结员	 果总数:2; j	<b>返回结果</b> 总	数:2;起始结果序号:	0						
doc_i	d 176									
id	211									
updat	etime 2012	2/09/12 0	0:00:00							
addre	制造 *SS 制造	商: 美国 商所属国家	北京中科泛华测控技术 <mark>恩艾仪器有限公司</mark> 邓: 美国 页目名称:电	<b>*有限公司</b>						
doc_i	d 177									
id	212									
updat	etime 2012	2/09/12 0	0:00:00							
addre	推荐 : · ·	:中标商: 页目名称:	北京中科泛华测控技才 B学量专用	*有限公司	制造商: 美	国恩艾仪器有阳	人可 制造商所属	属国家: 美国	招标内容	
1										



### 长句搜索

searchAgentClie		
服务器IP: 127.0	0 . 0 . 1 端口号: 8001 搜索編号: 0 毎页结果数:	20
字段文件:\IndexF	File\field.dat 显示格式:	html
	确认    默认值	
或法人	人委托书到汉寿县城关镇小南门55号(汉寿县建设局办公楼4楼)购买资格 搜索 下一页 上一页	退出
doc id	94	
id	118	
	e 2012/09/12 00:00:00	
address	r>常德工程招标代理有限公司受汉寿县土地开发整理中心委托,对其汉寿县围堤湖土地整理(二期工程)项目; 内公开招标。凡具有建设行政管理部门颁发的房屋建筑工程施工总承包叁级(00时(北京时间),逾期不到理。各投标人在网上报名后凭介绍信或法人委托书到汉寿县城关镇小南门55号(汉寿县建设局办公楼4楼)购买预	办
doc_id	95	
id	119	
updatetime	e 2012/09/12 00:00:00	
address	常德工程招标代理有限公司受汉寿县土地开发整理中心委托,对其汉寿县围堤湖土地整理(二期工程)项目进行公开招标。凡具有建设行政管理部门颁发的房屋建筑工程施工总承包叁级(00时(北京时间),逾期不予办各投标人在网上报名后凭介绍信或法人委托书到汉寿县城关镇小南门55号(汉寿县建设局办公楼4楼)购买资格	5理。

### 模糊搜索

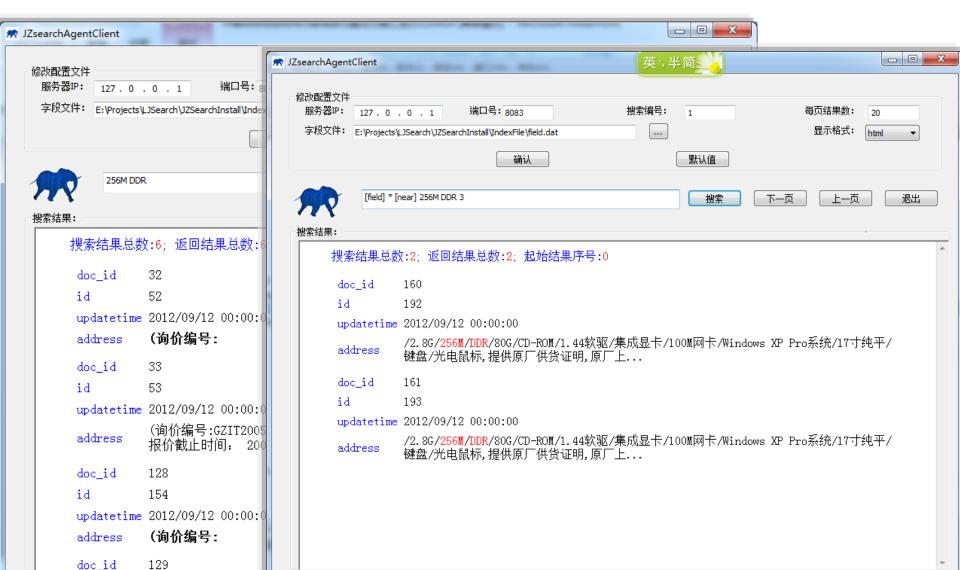
# 7公司名称表述不准确

searchAgent	tClient	中。半简单				
修改配置文件 服务器邛: 字段文件:	127 . 0 .	LJSearch\JZSearchInstall\IndexFile\field.dat 显示格式: html ▼				
		<b>确认</b> 默认值				
M	[field] * [f	fuzzy] 中设江苏机械设备集团公司				
搜索结果:						
搜索	*结果总数	文:173;返回结果总数:20;起始结果序号:0				
do	c_id	2				
id		3				
up	datetime	2012/09/12 00:00:00				
		<mark>中设江苏机械设备</mark> 进出口 <mark>集团公司</mark> 受买方委托对下列产品及服务进行国际公开竞争性招标。现 邀请合格投标人参加投标。				
ado	dress	1、招标产品的地点: 中 <mark>设江苏机械设备</mark> 进出口集团公司国际招标中心				
		5、投标截止时间和开标时间: 2005-07-2309:30:00				
		6、开标地点 <mark>:中设江苏机械设备</mark> 进				
do	c_id	3				
id		4				
up	datetime	2012/09/12 00:00:00				
ad	dress	中设江苏机械设备进出口集团公司受买方委托对下列产品及服务进行国际公开竞争性招标。现 邀请合格投标人参加投标。 1、招标产品的名称、地点:中设江苏机械设备进出口集团公 司国际招标中心 5、投标截止时间和开标时间: 2005-07-2309:30:00 6、开标地点:中设江苏 机械设备进				



# 设备参数精准搜索

7内存 256M的参数要求不等同于两者同时出现





# 语义关联搜索

### 7俗语洋蜡搜索到专业术语石蜡相关的内容



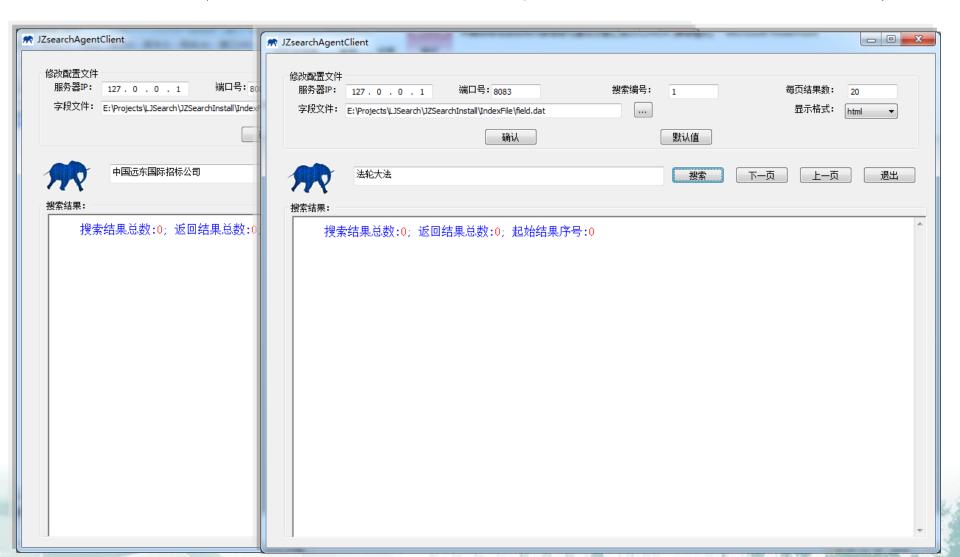
大数据

學LOGY



# 黑名单搜索屏蔽

7可以屏蔽禁止的公司、产品或反动信息搜索





# 正文、附件联合查询结果样例

标题:关于陕西某公司采购设备的公告

正文:

... 陕西法士特齿轮有限公司20万台变速器技术改造 3、数控高效流齿机 1台

附件:

065.txt

... 西安法士特汽车传动有限公司出口齿轮生产线技术改造项目 1、数控滚齿机11台

131.txt

... 陕西法士特齿轮有限公司20万台变速器技术改造 1、数控滚齿机 2台



### 搜索结果自动分组统计





加入收藏



设为首页







网站▼ | 登录 | 注册 | 返回首页

#### 标准分类

- 植物检疫、病虫害防治(18)
- 瓜果、蔬菜种植与产品(14)
- 农药管理与使用方法(12)
- 豆类、薯类作物与产品(10)
- 食品卫生(8)
- 农机具(7)
- 基础标准与通用方法(5)
- 农林技术(5)
- 种籽与首种(5)
- 粮食加工与制品(5)
- 经济作物综合(4)
- 畜禽饲料与添加剂(4)
- 植物保护综合(2)
- 一般有机化工原料(2)
- 基础标准与通用方法(2)
- 标志、包装、运输、贮存(2)
- 食品加工与制品综合(2)
- 标准化、质量管理(1)
- 粮食、饲料作物综合(1)
- 标准化、质量管理(1)
- 蔬菜加工与制品(1)
- 调味品(1)
- 蔬菜罐头(1)
- 烟草制品(1)

标准号	关键词: 土豆	检索 皮 按国内标	注 图 按国外标准
搜索结果自动分	分类检索 高級检索	首页 < 1 2 3	4   5   6   7   8   >   展示
统计 标准号		名称	中标分类名称
DB13/T 396.7-1999	旱地地膜覆盖栽培技术规程	自动搜索数据库里面的内	种籽与育种
DB13/T 398.7-1999	旱地地膜覆盖栽培技术规程	容,数据库更新10秒内,	农林技术
DB13/T 865-2007	马铃薯有机栽培技术规程	会直接反映在搜索结果	瓜果、蔬菜种植与产品
DB1302/T 156-2001	无公害马铃薯生产技术规程	里。	瓜果、蔬菜种植与产品
DB1304/T 071-2001	无公害农产品生产技术规程 马铃薯	标题	农林技术
DB1304/T 102-2001	马铃薯脱毒种薯繁育技术规程		农林技术
DB1304/T 136-2004	棉花与土豆(马铃薯)间套复种生产	技术规程	瓜果、蔬菜种植与产品
DB1306/T 11-2005	无公害马铃薯生产技术规程		瓜果、蔬菜种植与产品
DB1308/T 007-1999	马铃薯茎尖脱毒技术规程		豆类、薯类作物与产品
DB1308/T 008-1999	脱毒马铃薯微型种薯生产技术规程		豆类、薯类作物与产品
DB1308/T 009-1999	脱毒 <mark>马铃薯</mark> 基础种著生产技术规程		豆类、薯类作物与产品
DB1308/T 010-1999	脱毒 <mark>马铃薯</mark> 合格种薯生产技术规程		豆类、薯类作物与产品
DB1308/T 011-1999	马铃薯脱毒种薯质里标准		豆类、薯类作物与产品
DB1308/T 012-1999	马铃薯商品薯(块茎)质量标准		豆类、薯类作物与产品
DB1308/T 088-2005	食用马铃薯淀粉		豆类、薯类作物与产品

大数据分析与应用/张华平





### 我们的工作: JZSearch精准搜索引擎

	B政名址网 Database Center Of China Post		€ 返回主页	≥ 写信给我们 🗎	与我们联系
電 普通用户 300円户 高級用户	首页 邮政编码查询 企	事业单位查询 数据目录	数据信函业务	邮政黄页 学习园	地
用户名: 密码: 注册	<b>邮政编码查询 企事业单位查询</b> 中关村东路80号3号楼1509室	:	- <del>1</del>	<b>要素</b> 高級搜索	使用帮助
名址动态	● 邮 政 编 码 查 询			返回主页	查看数据
	□ 邮政编码查询结果列表:			查询出1条	符合条件的记录!
	<b></b>	<b>注荐结果:100190</b> 可能	能结果如下:		
※ 什么是邮政编码	邮政编码		行政区划		
MORE	100190	北京市海淀区中关村东路62~100	77		
新邮预订户	POWER BY AddreSmart	——— 第1页	[/共1页 共 <b>1</b> 条词		



### 我们的工作: JZSearch精准搜索引擎



● 全部 ○ TXT ○ DOC ○ PDF ○ HTML ○ PPT ○ Excel ○ 其他

全部共享

#### Unix vi.doc

...yw就是复制两个单词 如果要复制第m行到第n行之间的内容,可以在末行模式中输入m,ny例如:3, 🔒 5y复制第三行到第五行内容到缓存区。粘贴缓冲区中的内容,用p7.撤销操作 u命令取消最近一次的操作,可 以使用多次来恢复原有的操作 U取消所有操作 Ctrl+R可以恢复对使用u命令的操作 8.搜索及替换命令 vi的查找和替 换功能主要在末行模式完成: 至上而下的查找 / 要查找的字符雷,其中/代表从光标所在位置起开始查找,例 ? 要查找的字符雷 例如:/work 替换 如:/ work 至下而上的查找 :s/old/new用new替换行... JoinDOC\file\192.168.1.102\department\Unix vi... - 2010-11-09 14:44:49

来源:运维部-李新健 权限:所有

#### e 博士论文 杨少华.doc

... HYPERLINK \l "\_Toc228695035" - ¶图1.4 基于流程图建模的情景应用构造 ‼ PAGEREF \_Toc228695035 \h 🔒 ,¶21→→ ‼ HYPERLINK \l"\_Toc228695036",¶图1.5 情景应用构造系统的表达能力与易用性分析 ‼ PAGEREF Toc228695036 \h r¶33↓↓ !! HYPERLINK \l" Toc228695037" r¶图2.1 HTML信息源实例: Amazon图书 搜索 ! PAGEREF \_Toc228695037 \h [ 13944 !! HYPERLINK \l "\_Toc... JoinDOC\file\192,168,1.99\public\博士论文 杨少华,doc - 2010-11-09 14:55:45

来源:开发部:刘志华 权限:所有

#### 

...,挖掘网上论坛信息不仅可以便利广大网民搜索网络资源,同是对于政府掌握社情民意,构造和谐社会有 🚽 着重要的意义。与普通网站多以静态网页为主不同,论坛往往都是动态生成的。论坛的数据保存在后台数据 库中,根据用户提交的参数不同,动态从数据库中读取相关内容生成网页,因此论坛有其自己的特点。 1、链接层 次比较深。用户在浏览某一帖子时往往要点击数次才能找到所看的帖子,并且需要翻页数次才能浏览完毕,论坛采 集器也要模拟用户动作进行深层链接采集。 2、链接种类繁杂。论坛中除了有帖子对应的链接外,还存在着大量的 功能链接和其他链接等。如"评论"、"回复"等。 3、内容重复链接。论坛中往往存在大量链接不同而指向同

JoinDOC\file\192.168.1.99\public\论坛采集.doc - 2010-11-09 14:55:50

来源:开发部:刘志华 权限:所有

#### ▼ 文档统计

全部时间

PPT(4)	JPG(4)
GIF(19)	H(27)
HTM(28)	TXT(41)
DOC(51)	CPP(92)
PDF(163)	C(178)
BMP(254)	

找到相关结果约5篇,用时2.537秒

登录 🔥

- 量女童
- ▶ 部门统计
- ▶ 人员统计
- ▶ 文档TOP10
- ▶ 贡献TOP10
- 阅读TOP10

■ oracle管理.txt



北京理工大学 BEIJING INSTITUTE OF TECHNOLOGY



### 我们的工作: JZSearch精准搜索引擎

#### 少數民族语言發液技索引擎

偏移量	长度	文件路径	<b>存股</b>	全文
0	11532	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民网维吾尔文دونو نورى	(IEX)
57433	11532	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民网维吾尔文—خونى نورى	正文
69320	10674	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民阿維吾尔文 一 かりまり こうしゅう しょうしょう しょうしょう しょうしょ しょうしょ しょうしょ しょうしょ しょうしょ しょうしょう しょうしゃ しまいり しょうしょう しょうしょく しょうしょく しょうしょく しょくりょく しょくりょく しょくりょく しょくり しょくり しょくりょくりょくりょくりょくりょくりょくりょくりょくりょくりょくりょくりょくりょ	正文
80102	9776	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民阿维吾尔文-ئىسىدىمائىي خەۋەرلەرخەلقى تورف	正文
89986	10393	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民阿維吾尔文 しゅっかいきょう しゅっかい 人民阿維吾尔文	正文
101320	4984	F:/data/MultiLang/Corpus/20100426/000@2010042615	人民阿维吾尔文-مۇڧەت ئاسىراش-جەلق تورى	正文
106304	3690	F:/data/MultiLang/Corpus/20100426/000@2010042615	一人民阿維吾尔文 しゅっぺんしゅん	正文
109994	4304	F:/data/MultiLang/Corpus/20100426/000@2010042615	人民网维吾尔文-دوروس برافتادوني نوري	正文
114298	4037	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民网维吾尔文-شەخسەلەر-خەلق نورى	正文
118335	8526	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民阿维吾尔文-سۈرەنكىڭ خەۋەرلەرخەلق ئورى	正文
127751	4385	F:/data/MultiLang/Corpus/20100426/000@2010042619	人民阿·-دەنق تۈرگ-: amp;nbsp: سەشقىيە amp;nbsp; مەخسىۋىن	正文

下一页 尾页 共有 1380 条记录, 当前第1/92 页

باش بمت قىلىڭبمت ساقلاڭتەھرىر خىت ساندۇقىئالاقىلىشىڭباش بىتىمىركىز 人民网维吾尔文-خىلق تورى خەۋەرلىرىخىلقئارا خىۋەرلىرشىنجاڭ ;x0D; #x0A روھىمىملىكىت

خەۋەرلى رىئىقتىسادجەمئى يەتمائارى پتەنتەربى يەمىللەت ۋە مەدەنى يەتمۇ ھىت ئاسراشدى نقىزى ق ئۇقتاشە خسلەرسۈرەتلىك خەۋەرلەر مەلئى كەۋەرلەر ماتېرى يالمەخسۇس سەھى يەزەھەر ۋە ئەيدى زىپڭى خەۋەرلەر مەلئى كەتلىك ئەرلەر ۋاسكى تبول كەسپى يىرلەشمە مۇسابى قىسى ئەرلەر قاسكى تبول كەسپى يىرلەشمە مۇسابى قىسى ئەرلەر قاسكى تېرلەر قاسكى تېرلەشمە تەرپاندا كەم كۆرۈلى دىغان بورانلى قەۋارايى كۆرۈلدى گەنسۇ جاڭىبدا قۇم بورانلىق ھاۋارايى كۆرۈلدى 1-پەسىلە سىرتتىن شىنجاڭغا سېلىنغان مەلمەغ نوخشاش ۋاقىتتىكى دىن 4 يېرىم ھەسسە ئاشتى روسى يەنىڭ ئۇنجى تۈركۈمدىكى ياردەم بۇيۇملى رى يېتىپ كەلدى شىنجاڭ قۇمۇلدا قۇم بورانلىق ھاۋارايى كۆرۈلدى ئالدىن مەلۇمات: ئادىل ھوشۇر خەلق تورىدا ئىدىن دەلىق دەرىدا سەھدەت ئىلەپ بادى دۇرۇددى ئالدى مەلۇمات ئادىل ھوشۇر خەلق تورىدا ئىدىن دەلىق دۇرىدا ئىدىن مەلۇمات ئادىل ھوشۇر خەلق تورىدا



### 网上通用搜索功能

#### 2013北京旅游攻略 北京景点线路游记 百度旅游

最新<mark>北京</mark>旅游攻略,百度旅游为您准备了<mark>北京</mark>景点、线路、交通、美食、住宿等旅游攻略信息 和实用精美游记,供赴北京旅游的亲们参考。

lvvou.baidu.com/

#### 首都之窗-北京市政务门户网站

北海公园 图说<mark>北京 魅力公园 北京</mark>印象图集 "绿色出行"图片征集 园博会图片征集 2013 年<mark>北京</mark>冬季图片征集 热点关注[会议]市委十一届三次全会于12月22日至12月23...

www.beijing.gov.cn/ 2013-11-27 - V - 百度快照

#### 北京旅游攻略 北京北京旅游景点 北京旅游网

北京欣欣旅游网,提供北京北京旅游景点推荐、12月北京旅游攻略、北京旅行社、北京旅游线路、北京酒店预订、北京旅游地图等出行指南及旅游服务●欣欣旅游网 CNCN.com ... beijing.cncn.com/ 2013-12-07 ▼ ▼ - 百度快照

#### 北京汽车 北京车市 北京汽车报价 北京汽车网 汽车之家

汽车之家北京站为您提供北京汽车报价,北京汽车行情,北京汽车经销商推荐,最精彩的北京汽车新闻、评测、导购及二手车信息,是提供北京汽车信息最快的汽车网站www.autohome.com.cn/beijing/ 2013-12-07 ▼ - 百度快照

#### 2013北京旅游攻略,北京自助游攻略,蚂蜂窝北京出游攻略游记-蚂蜂窝

胡同是地道<mark>北京</mark>人的居所,这里充满着丰厚的生活气息,保留着许多文物古迹。当今的胡同,又进驻了许多潮流元素,这种传统与现代的结合,是<mark>北京</mark>最有魅力的地方之一。 .... www.mafengwo.cn/travel-scenic-spot/m... 2013-12-06 ▼ - 百度快照

#### 北京网-首都城市综合信息服务平台-北京生活、旅游、交通、文化、...

北京网是北京市政府主导建设,为公众提供首都城市综合信息服务的公益性网站。提供北京地图、交通、旅游、住房、餐饮、医疗、演出、教育、就业等领域的便民服务信息。

www.boijing.co/ 2013 12 07 - 古度出现





▲ 个人中心 🕊



综合搜索 组织结构 电网报 图片 电力地图 统计分析

刘振亚是谁

搜索

热搜词条:

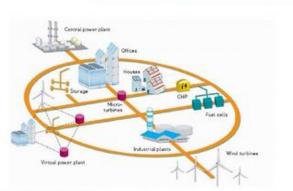
刘振亚是谁

UHV是什么

智能电网

配电网

全球能源互联网



#### 最新消息

国网大事记

- 严格依法治企 坚持以...
- 台风"灿鸿"肆虐供...
- 奋发有为 扎实工作 ...
- 连续抢修3天
- 风狂雨骤显担当
- 一定把损失隆到最低

北京理工大學 BEIJING INSTITUTE OF TECHNOLOGY



■聚类结果

特高压

教育实践活动

智能电网

\_

全球能源互联网

学习实践活动

经济社会发展

安全生产

#### ■ 语义结果

询问对象: 刘振亚,

刘振亚,男,汉族,1952年8月生,山东郯城人,1984年加入中国共产党,1971年参加工作,山东工学院电力系电力系统 及自动化专业毕业,大学学历,山东大学电气工程学院电气工程及其自动化专业硕士研究生毕业,电气工程硕士,教授级

高级工程师,享受国务院政府特殊津贴。现任国家电网公司董事长、党组书记。

#### ▶ 语义统计分析

分析对象: 刘振亚

#### 』 智能搜索

返回检索结果约1906个结果...

#### 1. 刘振亚总经理调研韩国STX集团大连造船厂项目的用电情况

来自栏目: 无设定栏目 板块: 要闻 发布时间: 2007/05/29 00:00:00 作者: 杜平

27日,<mark>刘振亚</mark>总经理(前排左四)一行在大连市市长夏德仁(前排左三)的陪同下,调研韩国STX集团大连造船厂项目的用电情况。

关键词: 刘振亚 命名实体--人物: 刘振亚#夏德仁#

#### 2. 刘振亚分别会见花旗 集团和通用电气高层

来自栏目: 无设定栏目 板块: 要闻 发布时间: 2009/04/28 00:00:00 作者: 岳文

#### □ 语义自动计算

相关新概念发现

沙捞

公司党组

电力集团公司

能源互联网

#### 相关人物聚类

习近平/61

李克强/22

舒印彪/17

李荣融/12

郑宝森/11

温家宝/10

俞正声/9

曹志安/9

帅军庆/7

高培/7

#### 相关作者聚类

姚雷/295

陶思遥/108

张超义/54

江莹/41

大数据分析与应用/张华平





● 语 V 绝计分析

▶ 净 语义统计分析

分割

ゲ 分析対象: 刘振亚

分析

全球能源互联网 换流阀

业务流程化 特高压电网

特高压

经济社会发展

特高压技术

电网发展

王学军

分析对象:刘振亚是谁

时间:2007-2015

分析机构:国家电网公司

数据来源:10年国家电网报

说明: 左图为关键词随着时间的 变化,关键词也在发生变化。





#### 3. 刘振亚会见马来西亚沙捞越州首席部长

**来自栏目:** 无设定栏目 板块: 要闻 发布时间: 2010/04/14 00:00:00 作者: 姚雷

总经理<mark>刘振亚</mark>在公司总部会见了到访的马来西亚沙捞越州首席部长泰益玛目一行,双方进行了亲切友好的交谈,并就发挥各自优势加强合作深入交换了意见。 刘振亚对马来西亚客人的来访表示热烈欢迎,并详细介绍了中国电力工业尤其是电网的发展情况。他希望双方能进一步增进了解和信任,达成合作共识,并在更广泛领域积极开...

关键词: 刘振亚 命名实体--人物: 刘振亚#哈吉·拉旺#杜至刚#

#### 4. 刘振亚会见东方电气集团董事长

来自栏目: 头条2 板块: 要闻 发布时间: 2009/07/16 00:00:00 作者: 姚雷

刘振亚对王计一行的到访表示热烈欢迎,对于东方电气广大干部员工在特大地震灾难面前顽强不屈的精神和恢复重建取得的成绩表示钦佩。刘振亚说,东方电气集团公司是我国最大的发电设备制造企业之一,国家电网公司一直十分关注东方电气集团灾后重建和发展,帮助东方电气集团恢复重建是国家电网的社会责任所在。经过大...

关键词: 刘振亚 命名实体--人物: 刘振亚#栾军#

#### 5. 刘振亚会见法国电力集团公司董事长

来自栏目: 无设定栏目 板块: 要闻 发布时间: 2010/12/09 00:00:00 作者: 姚雷

总经理<mark>刘振亚</mark>在公司总部会见了到访的法国电力集团公司董事长兼首席执行官亨利·普格里奥一行。双方表示,将巩固合作基础,加强交流,进一步拓展双方在电力领域的合作,实现未来共同进步。 刘振亚对亨利·普格里奥一行的来访表示热烈欢迎。他说,国家电网公司非常重视与法国电力同行的友好关系,双方已有的合作是双赢...

#### 相关概念词发现



#### 相关人物计算



基于语义的自动学习计算





### 感谢关注聆听!



张华平

Email: kevinzhang@bit.edu.cn

微博: @ICTCLAS张华平博士

实验室官网:

http://www.nlpir.org



大数据千人会