

# 信息检索与数据挖掘

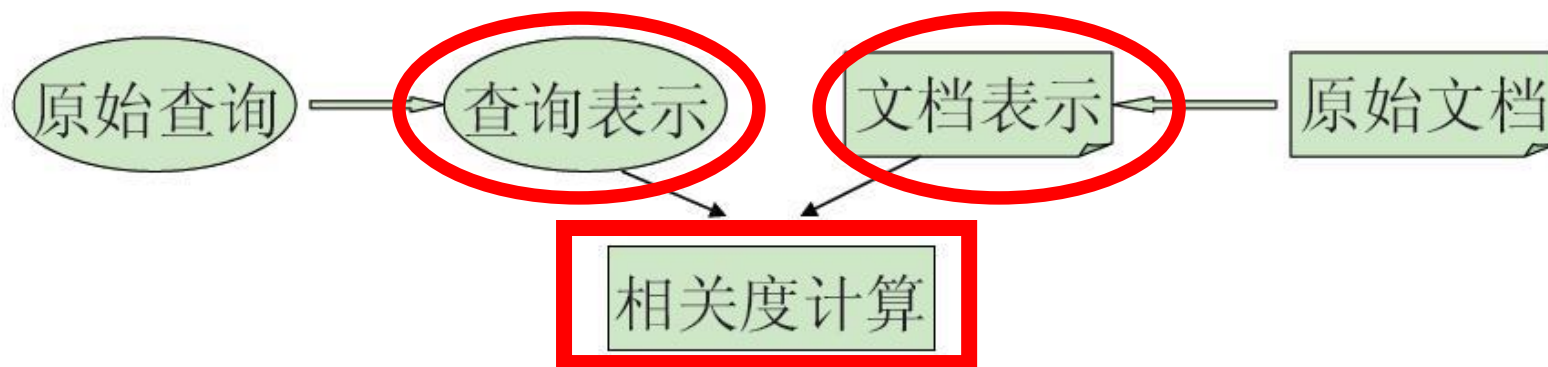
## 第3章 词项词典和倒排记录表

# 课程内容

- 第1章 绪论
- 第2章 布尔检索及倒排索引
- 第3章 词典查找及扩展的倒排索引
- 第4章 索引构建和索引压缩
- 第5章 向量模型及检索系统
- 第6章 检索的评价
- 第7章 相关反馈和查询扩展
- 第8章 概率模型
- 第9章 基于语言建模的检索模型
- 第10章 文本分类
- 第11章 文本聚类
- 第12章 Web搜索
- 第13章 多媒体信息检索
- 第14章 其他应用简介

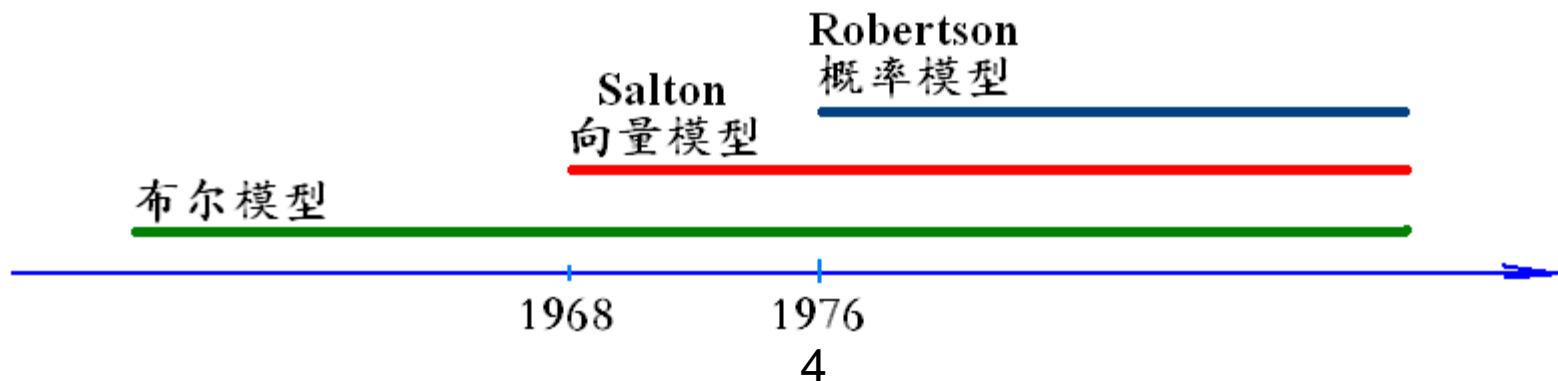
## 回顾2-1：信息检索模型的作用

- 信息检索模型是指如何对查询和文档进行表示，然后对它们进行相似度计算的框架和方法
  - 本质上是对相关度建模
- 信息检索模型是IR中的核心内容之一



## 回顾2-2：信息检索模型之经典模型

- 集合论模型 (Set Theoretic models)
  - 布尔模型 (Boolean Model, BM)、模糊集合模型、扩展布尔模型
- 代数模型 (Algebraic models)
  - 向量空间模型 (Vector Space Model, VSM)、广义向量空间模型、潜在语义标引模型、神经网络模型
- 概率模型 (Probabilistic models)
  - 经典概率论模型 (PM)、推理网络模型、置信网络模型



## 回顾2-3：线性扫描→词项文档索引

### Term-document incidence

- 解决方法是采用非线性的扫描方式，一种方法就是事先给文档建立索引（index）

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

如果文档（这里就是剧本）包含某个词，则对应的项为1，否则为0

***Brutus AND Caesar BUT NOT Calpurnia***

## 回顾2-4：词项文档索引→AND、OR、NOT

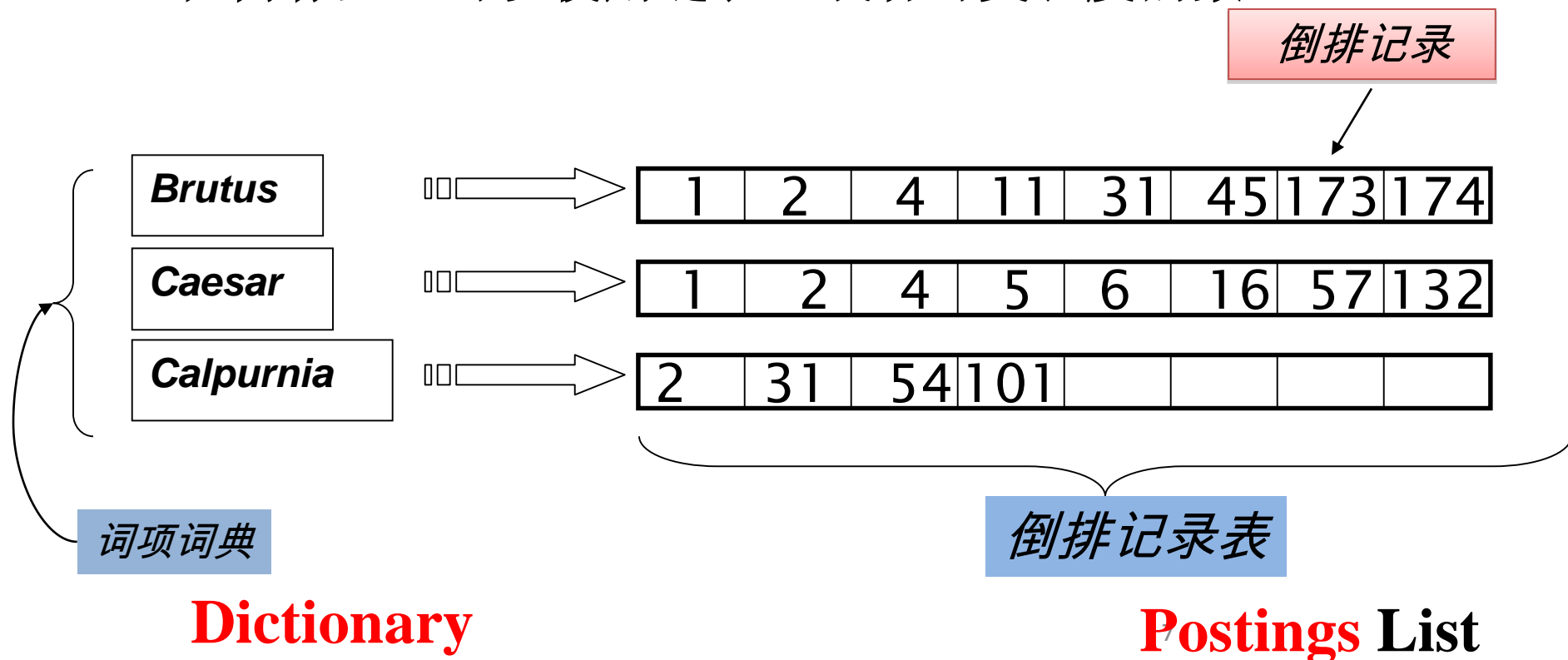
Term	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5	Doc 6	Doc 7	Doc 8
aid	0	0	0	1	0	0	0	1
all	0	1	0	1	0	1	0	0
back	1	0	1	0	0	0	1	0
brown	1	0	1	0	1	0	1	0
come	0	1	0	1	0	1	0	1
dog	0	0	1	0	1	0	0	0
fox	0	0	1	0	1	0	1	0
good	0	1	0	1	0	1	0	1
jump	0	0	1	0	0	0	0	0
lazy	1	0	1	0	1	0	1	0
men	0	1	0	1	0	0	0	1
now	0	1	0	0	0	1	0	1
over	1	0	1	0	1	0	1	1
party	0	0	0	0	0	1	0	1
quick	1	0	1	0	0	0	0	0
their	1	0	0	0	1	0	1	0
time	0	1	0	1	0	1	0	0

- **dog AND fox**
  - Doc 3, Doc 5
- **dog NOT fox**
  - Empty
- **fox NOT dog**
  - Doc 7
- **dog OR fox**
  - Doc 3, Doc 5, Doc 7
- **good AND party**
  - Doc 6, Doc 8
- **good AND party NOT over**
  - Doc 6

# 回顾2-5：超大词项文档矩阵→倒排索引

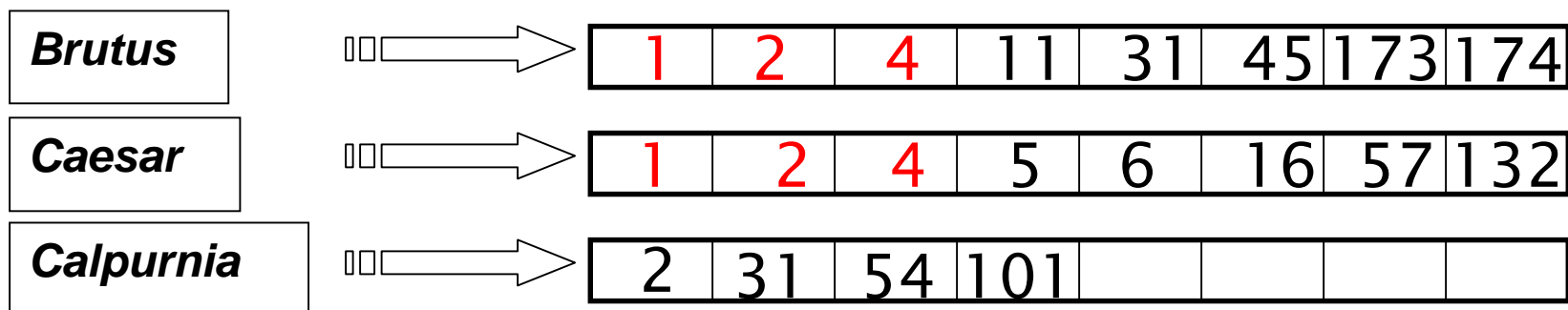
## Inverted index

- 应当使用可变长度的记录列表
  - 在硬盘上，一串连续的记录是正常的，也是最好的。
  - 在内存里，可以使用链表，或者可变长度的数组。



## 回顾2-6：如何基于倒排索引进行查询？→AND、OR、NOT

- Brutus AND Caesar



- Brutus OR Caesar 怎么求？

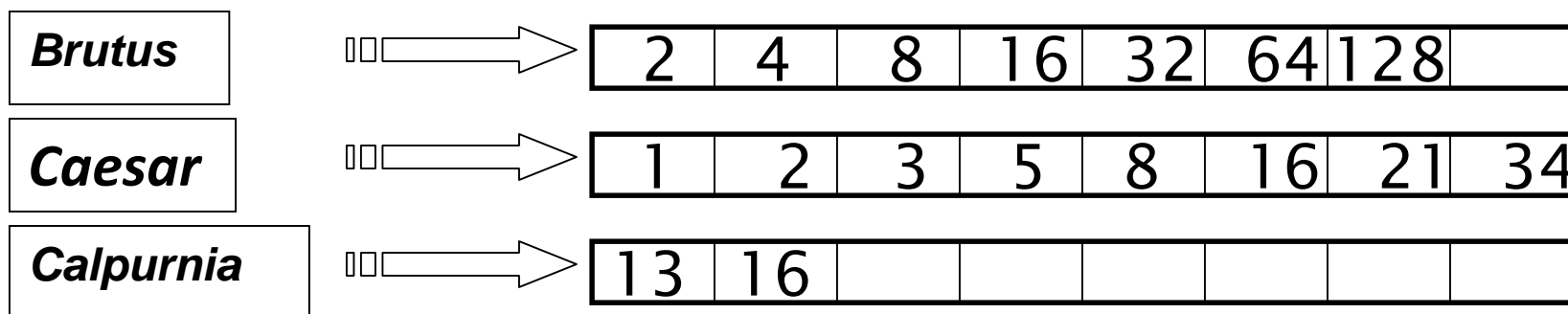
- NOT Brutus 怎么求？



## 回顾2-7：倒排记录表合并(merge)的优化→AND查询优化

- 按照文档频率的顺序进行处理：
  - 先处理文档频率小的，再处理大的

这就是为什么我们前面提到要  
存储词条的文档频率



按照(*Calpurnia AND Brutus*) *AND Caesar*的顺序处理查询

# Google中是否使用布尔模型？

- Google默认是与(AND)操作，输入查询 $[w_1 w_2 \dots w_n]$ 意味着  $w_1 \text{ AND } w_2 \text{ AND } \dots \text{ AND } w_n$
- 当返回文档不包含某个词 $w_i$  时，可能是如下情形：
  - 指向该页面的锚文本包含 $w_i$
  - 页面包含  $w_i$  的变形(不同形态的同一词，拼写校对，同义等等)
  - 长查询 ( $n$  large)
  - 布尔表达式返回的结果少
- 简单的布尔检索 vs. 结果的排序
  - 简单的布尔检索只返回匹配上的文档，不考虑结果顺序
  - Google和其他大部分精心设计的布尔引擎均对结果进行排序，以使好的结果排在差的结果的前面

# 上节课思考题

- 请推荐如下查询的处理次序

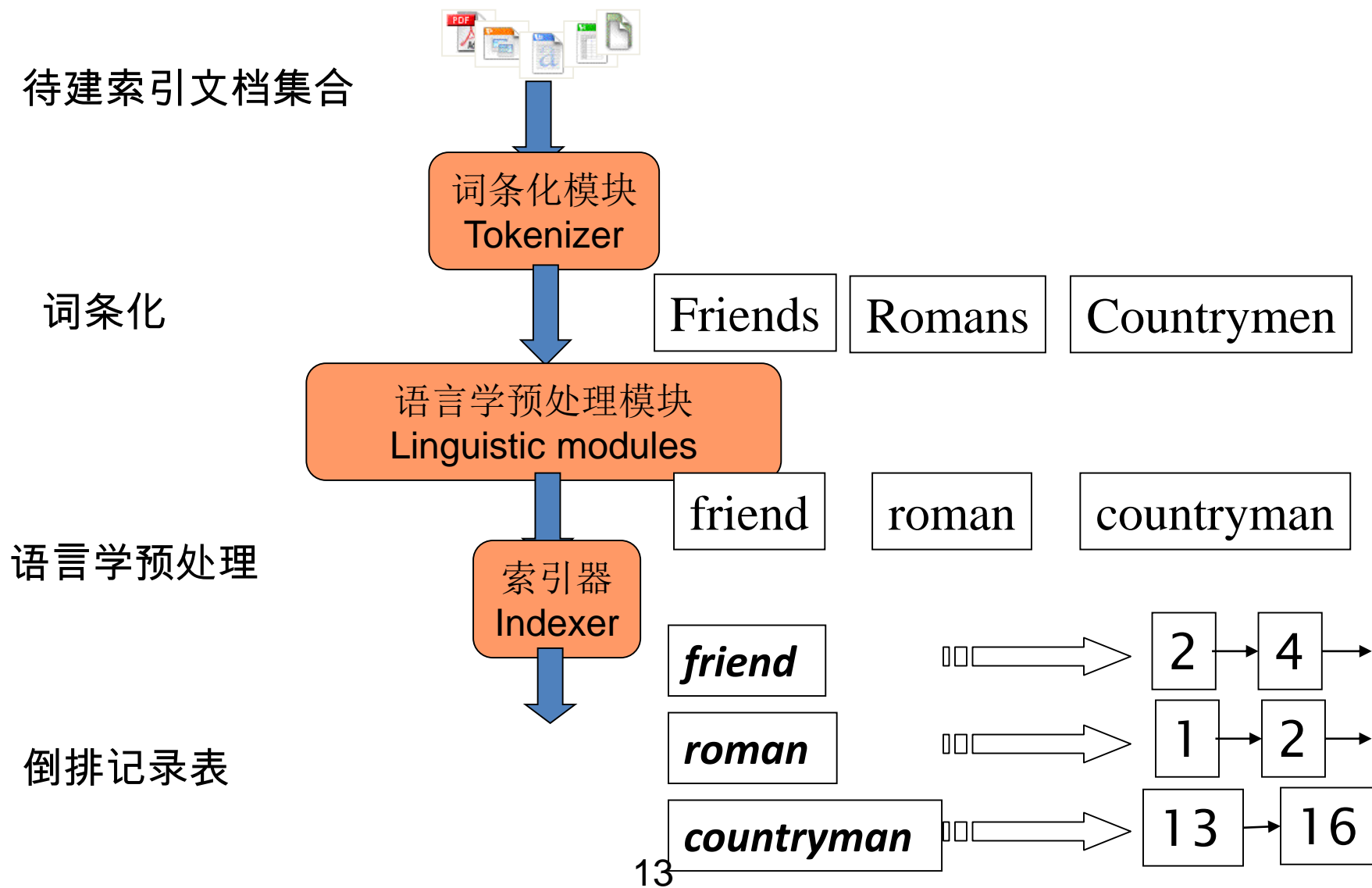
*(tangerine OR trees) AND (marmalade OR skies) AND (kaleidoscope OR eyes)*

Term	Freq
eyes	213312
kaleidoscope	87009
marmalade	107913
skies	271658
tangerine	46653
trees	316812

# 第3章 词典查找及扩展的倒排索引

1. 如何建立词项词典（ term vocabulary ）？
  - ① 文档集
  - ② 文本词条化（ Tokenization ）
  - ③ 语言学预处理
  - ④ 建立索引
  
2. 如何实现倒排记录表？
  - ① 快速合并算法：带跳表的倒排记录表(skip lists)
  - ② 包含位置信息的倒排记录表以及短语查询

# 建立词项 (Term) 词典过程



# 第3章 词典查找及扩展的倒排索引

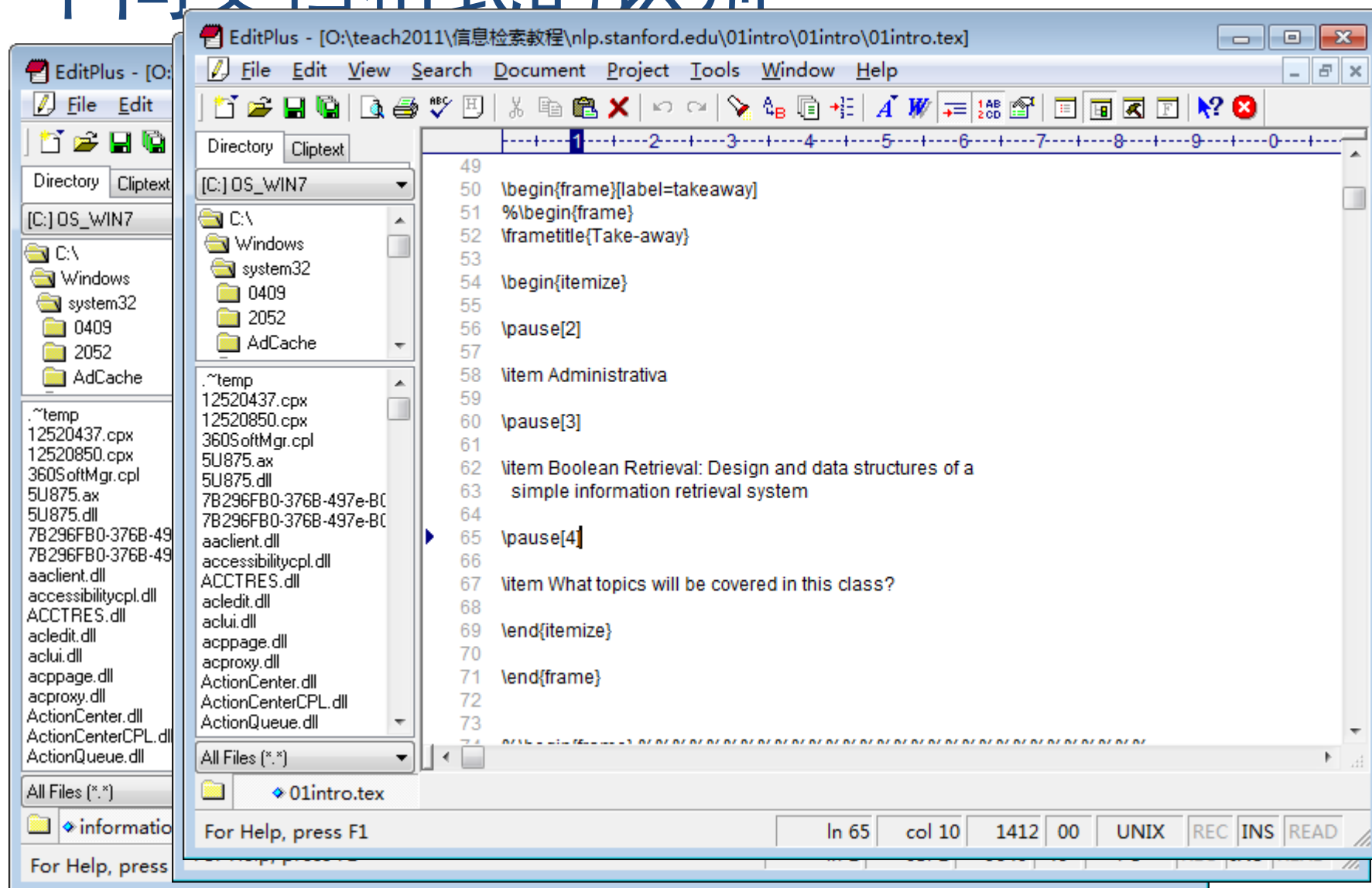
- 第一部分：如何建立词项词典？
  - 文档解析 (Parsing a document)
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并
  - 词干还原
- 第二部分：如何实现倒排记录表？
  - 快速合并算法：带跳表的倒排记录表
  - 包含位置信息的倒排记录表以及短语查询

# 文档解析

- 文档包含哪些格式？
  - pdf/word/excel/html?
- 文档中包含的语言？
- 文档使用何种编码方式？

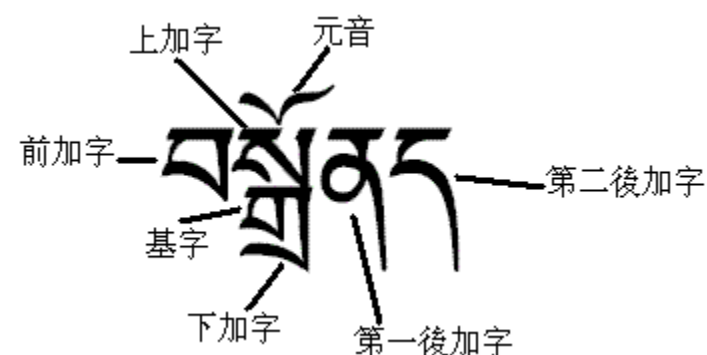
上述问题都可以看成是机器学习中的分类问题，但在实际中往往采用启发式方法来实现。（后面章节讨论）

# 不同文档格式的识别





# 文档中的语言



། ལྷོ་ཁྱིད་པོ་དུ་ཕེབས་རྒྱུ་དགའ་བསུ་ཞུ།

欢迎您到西藏来！

སྐྱ་ཁམས་བཟང་།

您好！早上好！下午好！晚上好！

བཀྲ་ཤིས་བདེ་ལེགས།

吉祥如意

དགོངས་པ་མ་ཚོ།

对不起

ཐུགས་རྗེ་ཆེ།

谢谢



# 文档中的编码方式

- 7bit ASCII?
- UNICODE?
  - UTF-8、UTF-16、UTF-32
- Email对二进制附件的编码
  - Content-Type: text/html;
  - charset="gb2312"
  - Content-Transfer-Encoding: base64

# 复杂因素：格式/语言

- 待索引文档集中包含不同语言的文档
  - 单独的一个索引应该包含不同语言的文档
- 一个文档或者其附件中包含多种语言或格式
  - 例子：一封法语的邮件中包含德语的pdf
- 文档单位的选择？
  - 一个文件？
  - 一封email？
  - 一封带有5个附件的email？
  - 一组文件？

# 第3章 词典查找及扩展的倒排索引

- 第一部分：如何建立词项词典？
  - 文档解析
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并
  - 词干还原
- 第二部分：如何实现倒排记录表？
  - 合并算法回顾
  - 基于跳表指针的快速合并算法
  - 短语查询

# 什么是词条化(Tokenization)

- 词条化：将给定的字符序列拆分成一系列子序列的过程，其中每一个子序列称之为一个“词条”。
- 输入：“*Friends, Romans and Countrymen*”
- 输出：
  - *Friends*
  - *Romans*
  - *Countrymen*
- 每个词条都作为候选的索引。
- 但是什么是有效的索引？

词条 (Tokens)  
词项 (Terms)

# 词条化可能遇到的问题 (英文)

e. g. : Finland' s capital →

Finland? Finlands? Finland' s?

- 连字符问题?

- Hewlett-Packard → Hewlett和Packard 是二个词条吗?
- State-of-the-art
- Co-education

- 空格问题?

- San Francisco是一个词条还是二个词条?

- 连字符和空格相互影响

- Lowercase, lower-case, lower case

- 英文句号的考虑

- IEEE 802.3 X.25 X.509

# 词条化可能遇到的问题(数字)

- 3/20/91    Mar. 12, 1991    20/3/91
- Tel:63601000            (800) 234-2333
- 查询2009至2011年间车祸死亡的人数
- B-52            AK-47
- PGP 密钥: 324a3df234cb23e
- 双11

# 词条化可能遇到的问题(中文)

- Out of Vocabulary
  - 人名、地名、机构名、一些新词
- Ambiguity
  - 同一句子有多种可能的分词结果
    - 南京市\_长江大桥 南京市长\_江大桥
    - 我们小组\_合成氢气 我们小\_组合成氢气
    - 发展中\_国家 发展\_中国\_家



# 词条化的策略

- 针对**不同的语言**，采取**不同策略**的词条化
- 分词的基本方法：
  - 基于词典的最大匹配法
  - 机器学习

正向最大匹配(基于词典的方法)

0 1 2 3 4 5 6  
他 说 的 确 实 在 理

逆向最大匹配(基于词典的方法)

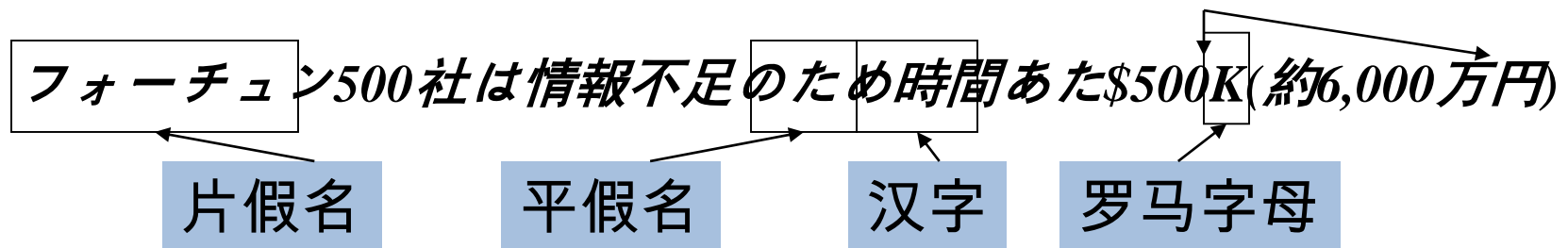
0 1 2 3 4 5 6  
他 说 的 确 实 在 理

指针位置	剩余词串	尾字	最大匹配词条
6	他说的确实在理	理	在理
4	他说的确实	实	确实
2	他说的	的	的
1	他说	说	说
0	他	他	他

指针位置	剩余词串	首字	最大匹配词条
0	他说的确实在理	他	他
1	说的确实在理	说	说
2	的确实在理	的	的确
4	实在理	实	实在
6	理	理	理

# 词条化可能遇到的问题(语言问题)

- 中文和日文词之间没有间隔：
  - 莎拉波娃现在居住在美国东南部的佛罗里达。
  - 分词结果无法保证百分百正确
- 日文中可以同时使用多种类型的字母表
  - 日期/数字可以采用不同的格式



而终端用户可能完全用平假名方式输入查询！

# 词条化可能遇到的问题 (语言问题)

- 阿拉伯文 (或希伯来文) 通常从右到左书写, 但是某些部分 (如数字) 是从左到右书写
- 词之间是分开的, 但是单词中的字母形式会构成复杂的连接方式

استقلت الجزائر في سنة 1962 بعد 132 عام من الاحتلال الفرنسي.

- 

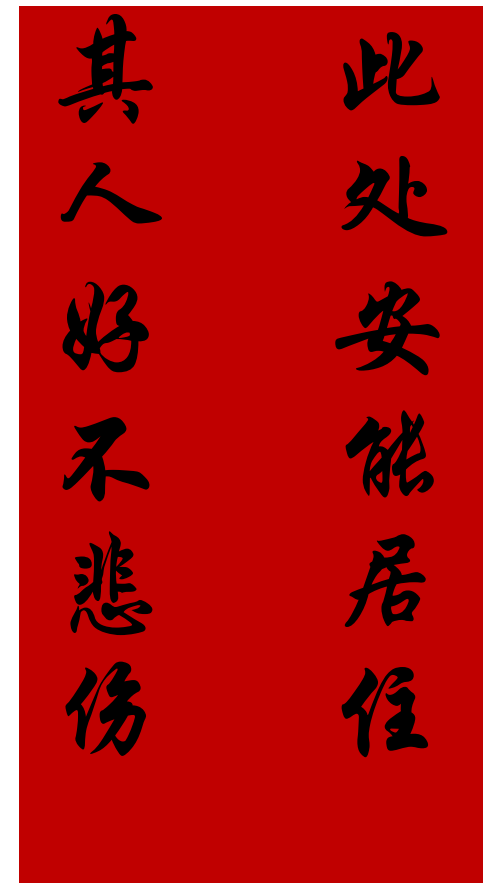
- ‘Algeria achieved its independence in 1962 after 132 years of French occupation.’

# 词条化可能遇到的问题(语言问题)

- 屈折语（俄语、英语、法语等）
  - 有比较丰富的词形变化；
  - 一种词形变化的语素可以表示几种不同的语法意义；
  - 词尾和词干或词根结合十分紧密
- 孤立语（以汉藏语系为代表）：
  - 词序严格；
  - 虚词重要；
  - 复合词多、派生词少
- 黏着语（日语、朝鲜语、蒙古语、维吾尔语等）
  - 只是词的尾部发生变化；
  - 一种变化只表示一种语法意义；
  - 词根与变词语素结合不很紧密，两者有很大的独立性

# 词条化不可协调的矛盾(二义性)

- 用红墨水写一个“蓝”字，
- 请问，这个字是红字还是蓝字？



# 第3章 词典查找及扩展的倒排索引

- 第一部分：如何建立词项词典？
  - 文档解析
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并
  - 词干还原
- 第二部分：如何实现倒排记录表？
  - 快速合并算法：带跳表的倒排记录表
  - 包含位置信息的倒排记录表以及短语查询

# 停用词

- 停用词表

- 将词项按照文档集频率 (collection frequency) , 从高到底排列
- 选取与文档意义不大, 高频出现的词, 比如, a, an, the, to, and, be...
- 停用词使用的趋势
- 现代搜索引擎发展的趋势使用少量的停用词表
- 现代IR系统更加关注利用语言的统计特性来处理常见词问题
  - 第五章介绍采用压缩技术降低停用词的存储开销
  - 第六章介绍词项权重, 将高频常用词来文档的排序影响降到最小
  - 第七章介绍一个索引按影响度大小排列的IR系统, 权重很小时, 停止扫描停用词表

# 消除停用词问题和可能的方法

- 优点：停用词消除可以减少term的个数
- 缺点：有时消除的停用词对检索是有意义的。
  - “的士”、“to be or not to be”
- 消除方法：
  - 查表法
  - 基于文档频率



# 第3章 词典查找及扩展的倒排索引

- 第一部分：如何建立词项词典？
  - 文档解析
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并
  - 词干还原
- 第二部分：如何实现倒排记录表？
  - 快速合并算法：带跳表的倒排记录表
  - 包含位置信息的倒排记录表以及短语查询

# 词项归一化

- 我们需要将文档和查询中的词条“归一化”成一致的形式
  - 希望USA和U. S. A. 之间也能形成匹配
  - 归一化的结果：
  - 在IR系统的词项词典中，形成多个近似词项的一个等价类
  - 隐式的建立等价类
    - 例如将USA和U. S. A. 映射为USA
    - 例如将anti-discrimination和antidiscrimination映射为antidiscrimination

# 词项归一化：不同语言之间的区别

- 重音符号：

e. g. : 法语中 *résumé* vs. *resume*.

- 变音符号：

e. g. : 德语中 *Tuebingen* vs. *Tübingen*.

（其实他们应该是等价的）

- 最重要的标准：

- 最重要的问题不是规范或者语言学的问题，而是用户将会如何根据这些词来构造查询？

- 即使是一些语言中，有的词有了标准的读音，但是用户有自己的读音/拼写方式

e. g. : *Tuebingen*, *Tübingen*, *Tubingen* → *Tubingen*

# 词项归一化：不同语言之间的区别

- 其他

- 中文中日期的表示7月30日 vs. 英文中7/30
- 日语中使用的假名汉字 vs. 中文中的汉字
- 词条化和归一化
- 二者都依赖于不同的语言种类，因此，在整个索引建立过程中要综合考虑
- e. g. :

*Morgen will ich in MIT...* 

**德语***Morgen will ich in MIT* 的意思是“我明天在MIT”，而德语中的“MIT”其实是“与”的意思

# 词项归一化：大小写转换

## ●一般策略

- 将所有字母转换为小写
- 绝大多数情况下，用户在构造查询时都忽略首字母的大写
- 一些专有名词除外
  - e. g. : *General Motors*
  - *Fed vs. fed*
  - *SAIL vs. sail*

## ●Google的例子

- 输入查询词C. A. T.
- 首页是关于猫的网站，而不是卡特彼勒公司（Caterpillar Inc.）（2005年的时候）



I keepz ur beerz till I getz toona

# 词项归一化

- 词项归一化的策略：建立同义词扩展表。

- 例子：

查询：

*window*

*windows*

*Windows*

检索：

*window, windows*

*Windows, windows, window*

*Windows*

# 扩展词表和soundex算法

- 如何处理同义词和同音词？

- e. g. : 手工建立同义词词表

- *car = automobile*                      *color = colour*

- ① 为每个查询维护一张包含多个词的查询扩展词表

- 例如：查询 *automobile* 的同时，也查询 *car*

- ② 在建立索引建构时就对词进行扩展

- 例如：对于包含 *automobile* 的文档，我们同时也使用 *car* 来索引

- 如何处理拼写错误？

- 其中的一种处理方法，就是根据发音相同来进行词项扩展
  - 后续章节中有讨论

# 第3章 词典查找及扩展的倒排索引

- 第一部分：如何建立词项词典？
  - 文档解析
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并 (Lemmatization)
  - 词干还原 (Stemming)
- 第二部分：如何实现倒排记录表？
  - 快速合并算法：带跳表的倒排记录表
  - 包含位置信息的倒排记录表以及短语查询



# 词形归并 (Lemmatization)

- 减少屈折变化的形式，将其转变为基本形式。
- e. g.
  - *am, are, is → be*
  - *car, cars, car's, cars' → car*
  - *the boy's cars are different colors → the boy car be different color*
- 词形归并可以减少词项词典中的词项数量

# 词干还原 (Stemming)

- 通常指的就粗略的去除单词两端词缀的启发式过程。
  - e. g. , *automate(s)*, *automatic*, *automation* → *automat*.

for **example** **compressed**  
and **compression** are both  
**accepted** as **equivalent** to  
**compress**.



for **exampl** **compress** and  
**compress** are both **accept**  
as **equival** to compress

# 中文重叠词还原——可视为“词干还原”

形容词(AB)	ABAB 式	AABB 式	A 里 AB 式
高兴	高兴高兴	高高兴兴	
明白	明白明白	明明白白	
热闹	热闹热闹	热热闹闹	
潇洒	潇洒潇洒	潇潇洒洒	
糊涂		糊糊涂涂	糊里糊涂
流气			流里流气
粘乎	粘乎粘乎	粘粘乎乎	
凉快	凉快凉快	凉凉快快	

形容词 (A)	ABB 式	ABCD 式
黑	黑压压	黑不溜秋
白	白花花	白不吡咧
红	红彤彤	
亮	亮晶晶	
恶	恶狠狠	
香	香喷喷	
滑	滑溜溜	

# Porter 算法

- 英文处理中最常用的词干还原算法。
  - 经过实践证明是高效性的算法。
- 算法包括5个按照顺序执行的词项约简步骤：
  - 每个步骤都是按照一定顺序执行的
  - 每个步骤中包含了选择不同的规则~~的约定~~
  - 比如，从规则组中选择作用时词缀最长的那条规则

C.J. van Rijsbergen, S.E. Robertson and M.F. Porter, 1980. New models in probabilistic information retrieval. London: British Library. (British Library Research and Development Report, no. 5587).

# Porter算法中的典型规则

- *sses* → *ss*      *care***sses** → *care***ss**
- *ies* → *i*      *pon***ies** → *po***in***i*
- *ational* → *ate*      *na***tional** → *na***te**
- 在这些规则中经常要考虑词的“测度”这一概念
- ( $m > 1$ ) *EMENT* →
  - *replacement* → *replac*
  - *cement* → *cement*

determi, determinate, determination, determinations, determine, determined, determines, determining → **determin**

support, supported, supportable, supportance, supporter, supporters, supporting, supportor, supports → **support**

# 一些词干还原工具

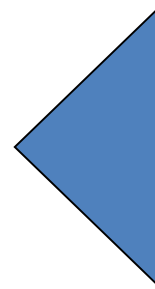
- 词干还原工具： Lovins
  - <http://www.comp.lancs.ac.uk/computing/research/stemming/general/lovins.htm>
  - 单遍扫描，最长后缀删除原则（大约包含250条规则）
  - 词干还原能够提高召回率，但是会降低准确率
  - e. g. :  
operative  $\Rightarrow$  oper
  - 词干还原对于芬兰语，西班牙语，德语，法语都有明显的作用，其中对芬兰语的提高达到30%（以MAP平均准确率来计算）。

# 语言的特殊性

- 词干还原和词形归并，都体现了不同语言之间的差异性，这些差异性包括：
  - 不同语言之间的差异
  - 特殊专业语言与一般语言的差异
  - 词干还原或者是词形归并往往通过在索引过程中增加插件程序的方式来实现
  - 商业软件
  - 开源软件

# Dictionary entries – first cut

<i>ensemble.french</i>
<i>時間.japanese</i>
<i>MIT.english</i>
<i>mit.german</i>
<i>guaranteed.english</i>
<i>entries.english</i>
<i>sometimes.english</i>
<i>tokenization.english</i>



These may be grouped  
by language (or not...).  
More on this in  
ranking/query processing.



# 本节内容小结： 如何建立词项词典？

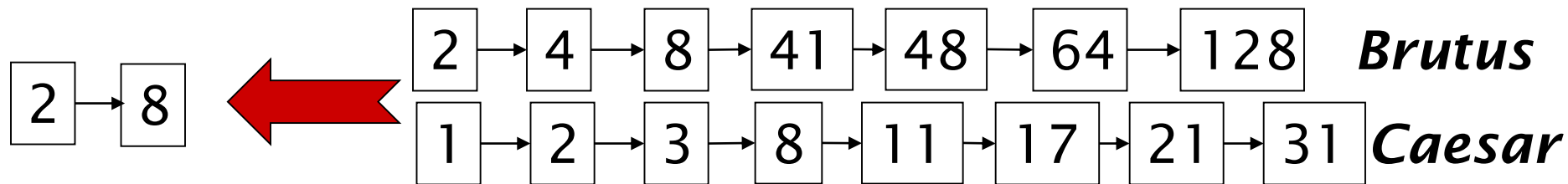
- 文档解析： 格式？ 语言？ 编码方式？
- 词条化
  - 概念词条 (Tokens) / 词项 (Terms)
  - 英文： 连字符？ 空格？ 句号？ 数字？
  - 中文： Out of Vocabulary? Ambiguity?
  - 方法： 针对不同的语言， 采取不同策略
- 停用词： 停用词表？ 查表法 or 基于文档频率
- 词项归一化：
  - 等价类？ 语言之间的区别？ 大小写转换？
  - 策略： 建立同义词扩展表
- 词形归并： am, are, is → be
- 词干还原： 去除单词两端词缀
  - Porter算法： 规则
  - 提高召回率， 但是会降低准确率

# 第3章 词典查找及扩展的倒排索引

- 第一部分：如何建立词项词典？
  - 文档解析
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并
  - 词干还原
- 第二部分：如何实现倒排记录表？
  - 快速合并算法：带跳表的倒排记录表
  - 包含位置信息的倒排记录表以及短语查询

# 合并算法 (Postings Merges) 回顾

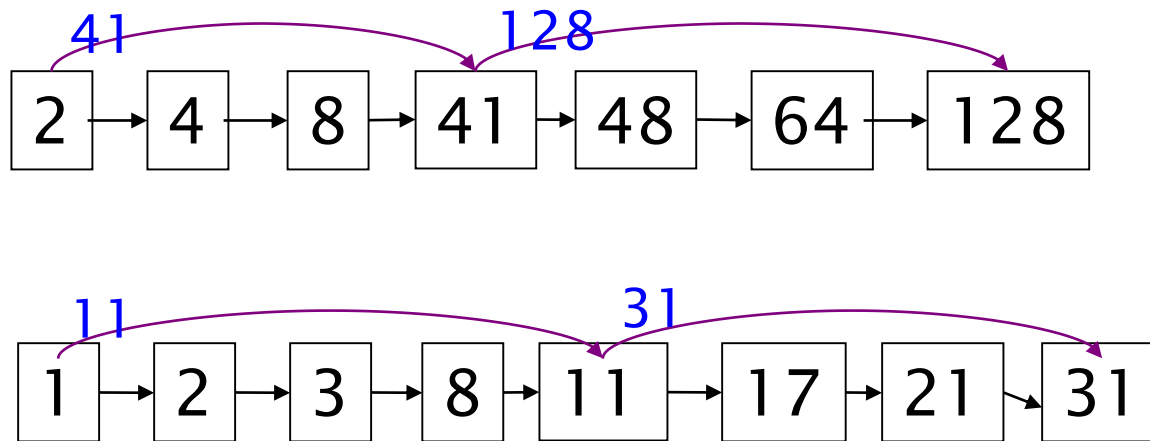
- 通过在二个倒排表之间同时移动指针来实现合并，此时的操作与线性表的总数成线性关系。



如果倒排表的长度分别是 $m$ 和 $n$ ，那么合并算法需要操作 $O(m+n)$ 次。

我们能否做的更好？ ← 上节课的问题

# 基于跳表（Skip List）的倒排记录表快速合并算法

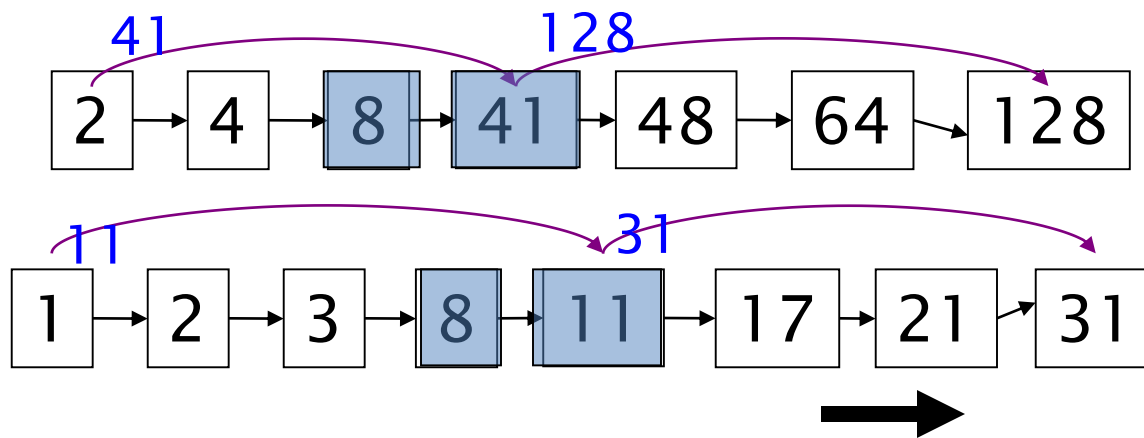


跳表指针能够跳过那些不可能出现在检索结果中的记录项。

构建跳表的二个主要问题：

- 如何利用跳表指针进行快速合并？
- 在什么位置设置跳表指针？

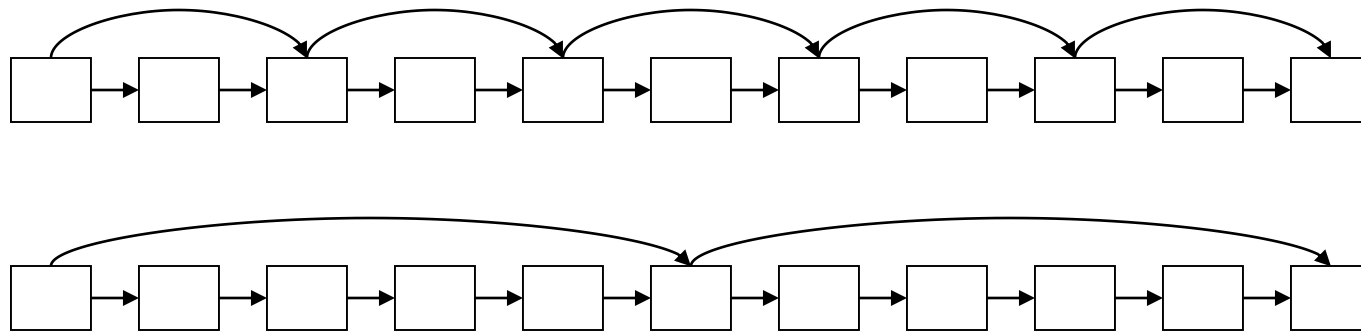
# 带有跳表指针的查询处理过程



1. 假定我们在进行遍历一直发现到了共同的记录8，将结果8放入结果表中之后，我们继续移动二个表的指针。
2. 假定第一个表指针移到41，第二个表的指针移到11。
3. 由于11比41小，因此，上面的表不需要继续移动，只需移动下面的表，跳到31。
4. 这样就跳过了17, 21。

# 在什么位置设置跳表指针？

- 策略：
  - 设置较多的指针→较短的步长⇒更多的跳跃机会
- 更多的指针比较次数和更多的存储空间
  - 设置较少的指针→较少的指针比较次数，但是需要设置较长的步长⇒较少的连续跳跃



# 设置跳表指针

- 放置跳表指针的一个简单的启发式策略是：

如果倒排表的长度是 $L$ ，那么在每个 $\sqrt{L}$ 处均匀放置跳表指针

- 该策略没有考虑到查询词项的分布
- 如果索引相对固定的话，建立有效的跳表指针比较容易，如果索引需要经常的更新，建立跳表指针就相对困难点。
- 硬件参数对索引构建有一定的影响
  - CPU速度
  - 磁盘访问速度

# 第3章 词典查找及扩展的倒排索引

- 第一部分：如何建立词项词典？
  - 文档解析
  - 词条化
  - 停用词
  - 词项归一化
  - 词形归并
  - 词干还原
- 第二部分：如何实现倒排记录表？
  - 快速合并算法：带跳表的倒排记录表
  - 包含位置信息的倒排记录表以及短语查询
    - 方法1：二元词索引
    - 方法2：位置信息索引



# 短语查询（Phrase Query）

- 用户希望将类似“stanford university” “中国科学技术大学”的查询中的二个词看成是一个整体。
- 类似“I want to university at stanford”这样的文档是不会被匹配的。
  - 大部分的搜索引擎都支持双引号的短语查询，这种语法很容易理解并被用户成功使用。
  - 有很多查询在输入时 没有加双引号，其实都是隐式的短语查询（如人名）。
  - 要支持短语查询，只记录 $\langle term : docs \rangle$  这样的条目是不能满足用户需要的。

# 第一种方法：二元词索引（ Biword indexes ）

- 将文档中每个连续词对看成一个短语
- 例如，文本“Friends, Romans, Countrymen”将生成如下的二元连续词对：
  - friends romans***
  - romans countrymen***
- 其中的每一个二元词对都将作为词典中的词项
- 经过上述的处理，此时可以处理二个词构成的短语查询

# 更长的短语查询

- 更长的短语查询可以分成多个短查询来处理
- 例如，文本 “***stanford university palo alto***” 将分解成如下的二元词对布尔查询：  
***stanford university AND university palo AND palo alto***
- 对于该布尔查询返回的文档，我们不能确定其中是否真正包含最原始的四词短语。



很难避免伪正例的出现！

# 扩展的二元词索引 (Extended Biword)

- 名词和名词短语构成的查询具有相当特殊的地位。
  - 首先对文本进行词条化，然后进行词性标注
  - 把每个词项分为名词 (N)、虚词 (X, 冠词和介词和其他词)。
  - 将形式为N\* $XN$ 非词项序列看成一个扩展的二元词
  - 每个这样的扩展的二元次对应一个词项
- 例如: *catcher in the rye*(书名: 麦田守望者)  
$$N \quad X \quad X \quad N$$
- 利用这样的扩展二元词索引处理查询,
  - 将查询拆分成N和X
  - 将查询划分成扩展的二元词
  - 最后在索引中进行查找

## 二元词索引的问题：

- 会出现伪正例子
- 增加词汇表的大小
  - 由于词典中词项数目剧增，导致索引空间也激增
  - 如果3词索引，那么更是空间巨大，无法忍受
- 二元词索引并非标准的解决方案，后面讨论的复合索引机制可以更加完美的解决短语查询的问题。

## 第二种方法：位置信息索引（Positional indexes）

- 在这种索引中，对每个词项，采取以下方式存储倒排表记录：

<词项，词项频率；

文档1：位置1，位置2……；

文档2：位置1，位置2……；

……>

<*be*: 993427;

*1*: 7, 18, 33, 72, 86, 231;

*2*: 3, 149;

*4*: 17, 191, 291, 430, 434;

*5*: 363, 367, ...>

# 位置信息索引例子

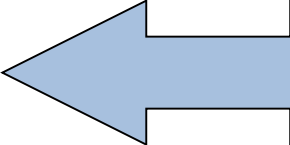
<*be*: 993427;

*1*: 7, 18, 33, 72, 86, 231;

*2*: 3, 149;

*4*: 17, 191, 291, 430, 434;

*5*: 363, 367, ...>



Which of docs *1,2,4,5*  
could contain “*to be*  
*or not to be*”?

- 对于短语查询，仍然采用合并算法，查找符合的文档
- 不只是简单的判断二个词是否出现在同一文档中，还需要检查他们出现的位置情况

# 短语查询过程

- 例子:

查询词: “to be or not to be”

倒排表:

- *to*:

2:1, 17, 74, 222, 551;

4:8, 16, 190, 429, 433;

7:13, 23, 191; ...

- *be*:

1:17, 19;

4:17, 191, 291, 430, 434;

5:14, 19, 101; ...

1. 考虑to和be的倒排表的合并, 查找同时包含to和be的文档
2. 检查表中, 看看是否某个be的前面的一个位置上正好出现to



# 邻近查询 (Proximity queries)

- Employ me/3 place, 表示从左边或右边相距在k个词之类
- 显然，位置索引能够用于邻近搜索，而二元词搜索则不能
- 临近查询在上一节课中有示例

# 回顾示例：WestLaw



<http://www.westlaw.com/>

- 最大的收费的法律搜索服务提供商
- 几十T的数据；700,000多用户
- 大多数用户仍然使用布尔查询
- 查询的例子：
  - What is the statute of limitations in cases involving the federal tort claims act?
  - **LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM**
    - **/3** = within 3 words, **/S** = in same sentence
- 许多专业的用户通常喜欢使用布尔查询
  - 因为用户精确的知道自己会得到什么
- 但是这并不意味着能得到更好的结果

# 位置信息索引的讨论

- 采用位置索引会大大增加倒排记录表的存储空间，即使采用后面讨论的压缩方法也无济于事。
- 由于用户期望能够进行短语查询和邻近查询，所以还是得采取这种索引方式。
- 另外，位置索引目前是实际检索系统的标配，这是因为实际中需要处理短语(显式和隐式)和邻近式查询

# 位置信息索引的大小

- 位置索引需要对词项的每次出现保留一个条目，因此索引的大小取决于文档的平均长度。
- 网页的平均长度不超过1000个词项。
- 但是某些文件，（SEC 股票文件）很容易就达到1000,000个词项。
- 假设，平均1000个词项中每个词项的频率都是1

Document size	Postings	Positional postings
1000	1	1
100,000	1	100

索引空间相差两个数量级。 69

# 经验法则(English-like)

- 位置索引大概是非位置索引大小的2—4倍
- 位置索引的大小大约是原始文档的30%—50%
- 提醒：上述经验规律适用于英语及类英语的语言

# 混合索引机制

- 二元词索引和位置索引二种策略可以进行有效的合并
  - 对于高频查询词可以采用二元次索引，例如“***Michael Jackson***”，
  - Williams等人（2004）评估了更复杂的混合索引机制，（引入后续词索引方法）。
  - 对于一个典型的web短语混合查询，其完成时间大概是只使用位置索引的1/4
  - 比只使用位置索引增加26%的空间

# 本节内容小结

- 带跳表的倒排记录表
  - 跳表 (Skip List)
  - 跳表指针：位置？个数？
  - 如果索引需要经常的更新？
- 包含位置信息的倒排记录表
  - 短语查询→二元词索引
    - 二元词索引→扩展的二元词索引：词性标注
    - 增加词汇表的大小
  - 短语查询→位置信息索引
    - 位置信息索引→邻近查询
    - 大大增加倒排记录表
  - 短语查询→混合索引机制

# 本章内容小结

- 如何建立词项词典？
  - 文档解析：格式？语言？编码方式？
  - 词条化：词条 (Tokens) / 词项 (Terms)
  - 停用词：停用词表？查表法 or 基于文档频率
  - 词项归一化：等价类  $\leftrightarrow$  同义词扩展表
  - 词形归并：am, are, is  $\rightarrow$  be
  - 词干还原：去除单词两端词缀、Porter 算法
- 如何实现倒排记录表？
  - 跳表：跳表指针 (位置、个数、更新问题)
  - 短语查询
    - 二元词索引  $\rightarrow$  扩展的二元词索引：词性标注
    - 位置信息索引  $\rightarrow$  邻近查询
    - 增加倒排记录表
    - $\rightarrow$  混合索引机制