

计算机应用编程实验一

大规模字符串



熊永平@网络技术研究院

ypxiong@bupt.edu.cn

周四15:30-17:30@3-117

2016.9.22

课程方式

- ▶ 总人数49人
 - □ 分为7组
 - □ 每组7人,讨论技术方案和调研
 - □ 每组分为若干实验单元,每单元2-3人,评分单元
- > 讨论课
 - □ 每组派人做报告
 - □ 由组长安排
 - □ 组长名单: 待定
- > 助教
 - □ 分发课程材料

课程表

学 期						क्रो	7人_			1		· [] [] -	_	期		क्रे	7人 -	<u>.</u>			勺	マゴム	. ПП		F 1	叚	
年 份						大	验-			L	头	验_			L	头	验:) •	字验 ——			七	年	
月份	八月		7	九月	1			+	月			+ -	一月	I		F	=	月				月			二	月	
周次	0	1	11	11	四	Į.	六	七	八	九	+	+ -	+=	+ 111	1	十五	十六	+4	十八	1/	=+	‡ 1	‡ 11	# =	廿四	廿五	廿六
星期一	22	本科生	研究生 上课	12	本科: 上記	26		10	17	24	31	7	14	21	8	5	12	19	26	/2	9	16	23	30	6	13	20
星期二	23	新生报到	6	13	20	27		11	18	25	1	8	15	22	29	6	13	20	1/	3	10	17	24	31	7	14	21
星期三	24	本科生 开学典礼	7	14	2	28	国庆节 假期	12	19	26	2	9	16	23	30	7	14	21	28	4	11	18	25	1	8	15	22
星期四	25	研究生 新生报到	8		22	29		13	20	27	3	10	17	24	1	8	15	22	29	5	12	19	26	2	9	16	23
星期五	26	研究生 开学典礼		中秋节 假期	23	30		1	21	28	4	1	18	25	2	9	16	23	30	6	3	20	27	3	10	17	24
星期六	27	3	þ		24	国庆节	8	1	22	29	5	12	19	26	3	10	1	24	31	7	14	21	春节	4	11	18	25
星期日	28	4		18	25	假期	9		23	30	6	13	b	27	4	11	18	5	元旦	8	15	22	29	5	12	19	26

课程介绍

实验一讨论课

实验二讨论课

实验三讨论课

实验四讨 论结束

实验一:海量字符串查找



实验背景

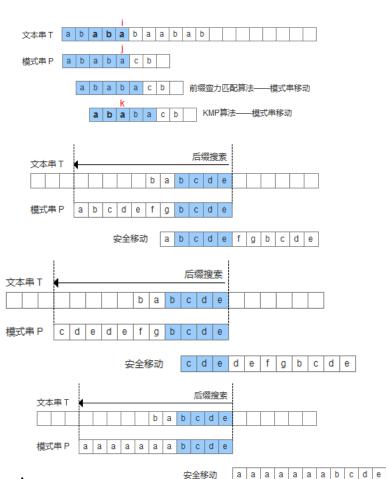
- > 一些实际问题
 - □ 如何实现在搜索引擎输入框自动提示?
 - □ 一个查询串的重复度越高,说明查询它的用户越多,也就越热门。大型搜索引擎每小时有几十亿个查询请求,如何统计搜索引擎最热门前100个查询?
 - □ 在实现一个编辑器时,如何对输入的单词进行拼写检查?
 - □ GFW每次google中输入**词后,如何在**名单中查找并reset?
 - □ 搜索引擎的网络爬虫,每天要爬取几十亿网页,哪些URL是爬过的?
 - 收到一封邮件后,能否快速在几亿个垃圾邮件黑名单地址里快速判断发件人是否在黑名单里面?
 - □ 检测引擎中包含几千万条特征字符串规则,如何在10G网络流量环境下检测网络流中的恶意软件特征?

实验一:海量字符串查找

- ▶ 问题
 - □ 在给定的海量个数的字符串中查找特定的字符串
- ▶ 挑战
 - □ 实际需求
 - 几亿规模
 - □ 数据量大
 - 200,000,000量级
 - □ 外存便宜
 - 存储成本低
 - □ 内存不够大?
 - 200,000,000*40bytes=8000,000,000bytes=8G

根据以往的学习思路

- > 字符串集合转换成一个大字符串
- > 字符串匹配
 - Brute Force
 - Strstr()
 - KMP算法
 - 前缀匹配算法
 - Boyer-Moore
 - 后缀匹配算法
- ▶ 问题
 - □ 存储空间
- > 复杂度
 - best case O(n/m)
 - worst case still O(nm)(BM) \ O(n)(KMP)





2、核心算法之Trie树

3、核心算法之BloomFilter

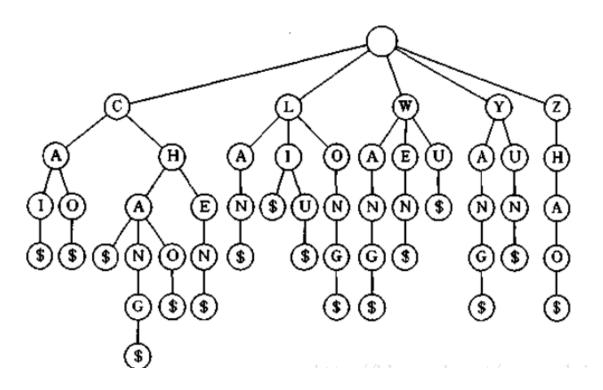
4、代码提交说明

Trie树

- ▶ 基于关键码分解的数据结构,叫作Trie结构(Trie树)
- ▶ 基于两个原则
 - 有一个固定的关键码集合
 - □ 对于结点的分层标记
- > Trie树
 - □ 又称单词查找树、字典树,是一种树形结构,是一种用于快速 检索的多叉树结构
- ▶ 典型应用
 - □ 统计和排序大量的字符串
 - 文本词频统计和文本检索
 - □ 优点:最大限度地减少无谓的字符串比较,查询效率比哈希表高。

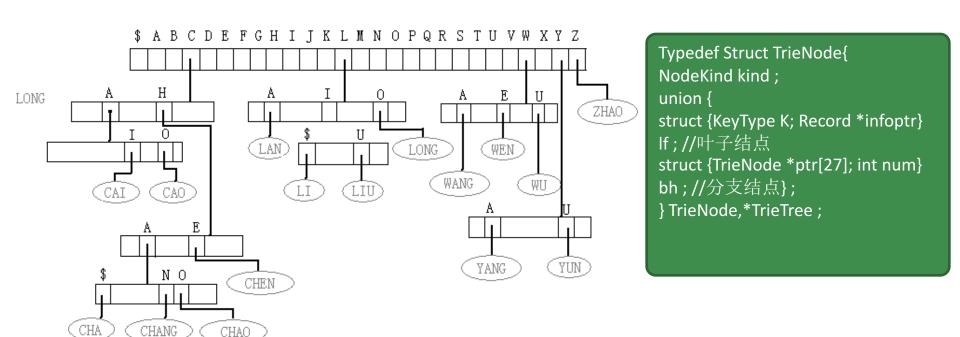
Trie结构示例

- > 存储词典
 - ☐ { CAI、CAO、LI、LAN、CHA、CHANG、WEN、CHAO、YUN、YANG、LONG、WANG、ZHAO、LIU、WU、CHEN }
 - □ 树的高度为最长字符串长度



Trie的实现

- 数据结构实现
 - □ 分支结点: 含有d个指针域和一个指示该结点中非空指针域的个数的整数域。
 - 分支结点所表示的字符是由其指向子树指针的索引位置决定的叶子结点:含有关键字域和指向记录的指针域。



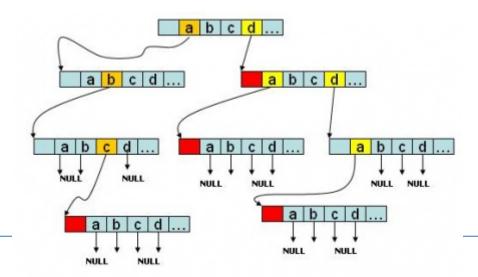
介绍

- > 不依赖于关键码的插入顺序
 - □ 树的深度受到关键码精度的影响
- ▶ 最坏的情况下,深度等于存储关键码所需要的位数
 - □ 例如,如果关键码是0到255之间的整数,关键码的精度就是8 个二进制位。
 - □ 如果有两个关键码: 10000010和10000011,它们的前面7位都 是相同的
 - □ 所以直到第8次划分才能将这两个关键码分开
 - □ 这样的搜索树深度也为8,但这是最坏的情况
 - □ 与B+树一样,基于关键码空间分解的树结构,其内部结点仅作 为占位符引导检索过程,数据纪录只存储在叶结点中

Trie查找

▶ 查找

- □ 在Trie树上进行检索总是始于根结点。
- □ 取得要查找关键词的第一个字母,并根据该字母选择对应的子树并转 到该子树继续进行检索。
- 在某个结点处相应的子树上,取得要查找关键词的第二个字母,并进一步选择对应的子树进行检索。
- □ 键词的所有字母已被取出,则读取附在该结点上的信息,即完成查找
- ▶ 示例
 - □ trie树中保存了abc、d、da、dda四个单词



Trie树的插入

> 插入

- 首先根据插入纪录的关键码找到需要插入的结点位置
- □ 如果该结点是叶结点,那么就将为其分裂出两个子结点,分别 存储这个纪录和以前的那个纪录
- 如果是内部结点,则在那个分支上应该是空的,所以直接为该分支建立一个新的叶结点即可

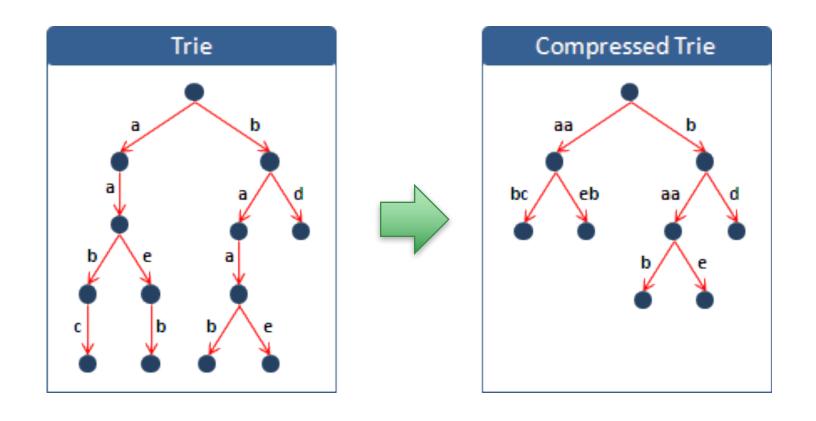
Trie查找效率分析

- ➤ 在trie树中查找一个关键字的时间和树中包含的结点数无关,而取决于组成关键字的字符数。 (对比:二叉查找树的查找时间和树中的结点数有关O(log₂ⁿ)。)
- ➤ 如果要查找的关键字可以分解成字符序列且不 是很长,利用trie树查找速度优于二叉查找树。
- ➤ 若关键字长度最大是5,则利用trie树,利用5次比较可以从26⁵=11881376个可能的关键字中检索出指定的关键字。而利用二叉查找树至少要进行log₂26⁵=23.5次比较。

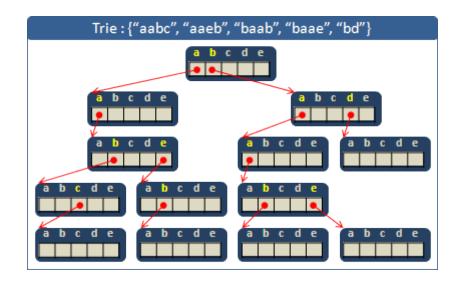
Trie树特性

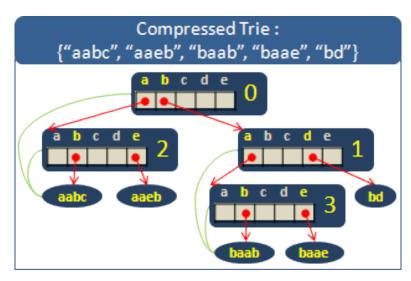
- ▶ 核心思想
 - □ 空间换时间
 - 利用字符串的公共前缀来降低查询时间的开销以达到提高效率的目的
- ▶ 优点
 - □ 查找效率高,与词表长度无关
 - Trie树的查找效率只与关键词长度有关
 - □ 索引的插入,合并速度快
 - 注意,直接遍历Trie树需要搜索大量的无效节点
 - 可以把数据存在一个数组中, Trie只保存指针
 - 这样合并时,只需要对数组进行遍历即可
- > 缺点
 - □ 内存空间消耗大
 - 如果是完全m叉树,节点数指数级增长
 - 不可达上限: 词数×字符序列长度×字符集大小×指针长度
 - 例如:20000×6×256×4=120M
 - □ 实现较复杂

Trie优化(1)--Compressed Trie



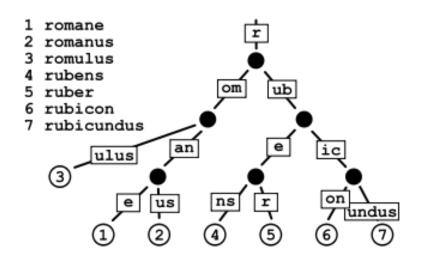
Radix Tree

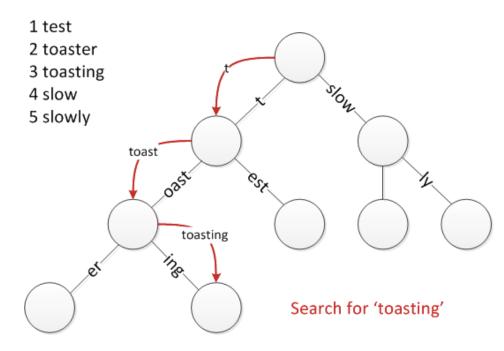




- 去掉没有分岔、连成直线的节点。
- 每个节点增加一个数字,纪录是第几个字符开始分岔。
- 去掉节点之后,字符串信息不完整,在树叶里储存完整字符串。
- 每个节点增加一个指针,纪录要参考哪一个叶子的字符串开头。

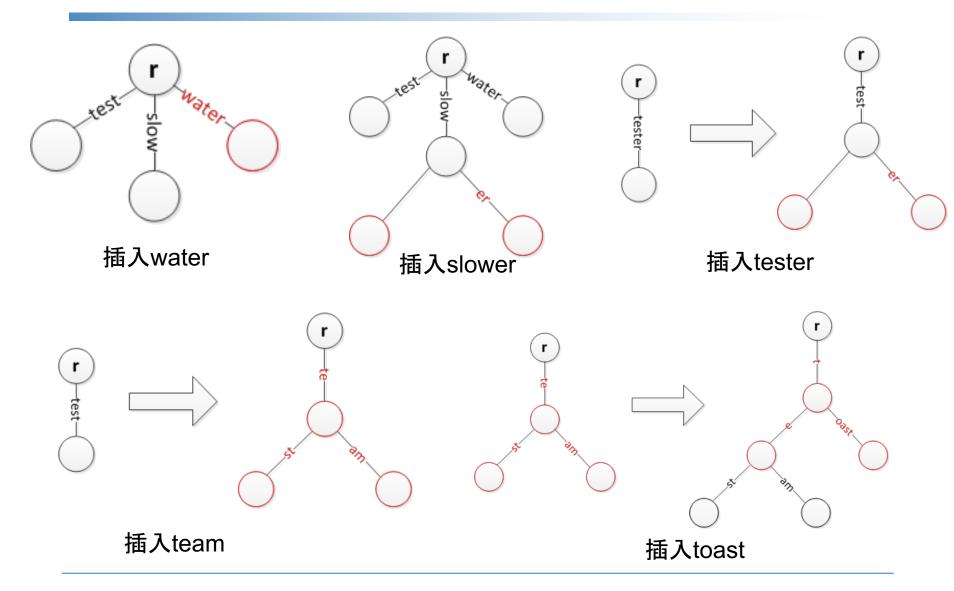
Radix Tree





Radix Trie 查找

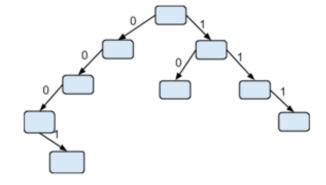
Radix Tree插入



Trie优化(2)-PAT Trie

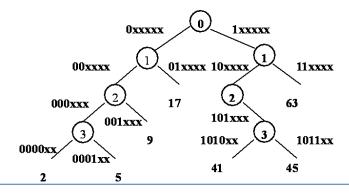
- Trie结构缺点
 - □ Trie结构显然也不是平衡的
 - □ 存取英文单词时,显然t子树下的分支比z子树下的分支多很多
 - □ 26个分支因子使得树的结构过于庞大,检索不便
- PATRICIA Trie ("Practical Algorithm To Retrieve Information Coded In Alphanumeric")
 - 关键码二进制形式存储
 - □ 根据关键码每个二进制位的编码来划分
 - □ 是对整个关键码大小范围的划分

每个内部结点都代表一个位的比较,必然产生两个子结点,所以它是一个满二 叉树,进行一次检索,最多只需要关键码位数次的比较即可。



PATRICIA应用举例

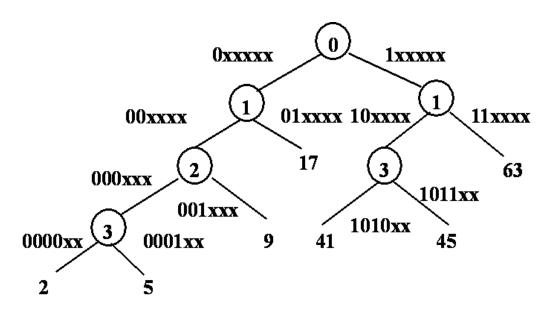
- ▶ 举例(2、5、9、17、41、45、63)
 - □ 因为最大的数是63,用6位二进制表示即可
 - □ 每个结点都有一个标号,表示它是比较第几位,然后根据那一位是0还是1来划分左右两个子树
 - □ 标号为2的结点的右子树一定是编码形式为xx1xxx, (x表示该位或0或1, 标号为2说明比较第2位)
 - □ 在图中检索5的话,5的编码为000101
 - □ 首先我们比较第0位,从而进入左子树,然后在第1位仍然是0 ,还是进入左子树,在第2位还是0,仍进入左子树,第3位变 成了1,从而进入右子树,就找到了位于叶结点的数字5



PATRICIA压缩优化

▶ 优化

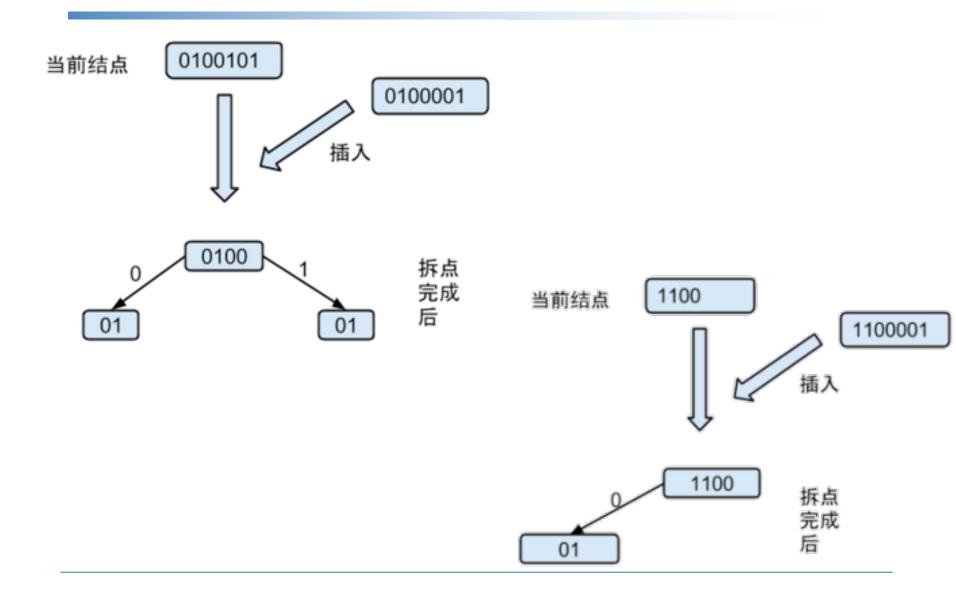
□ 在区分2和5、41和45时,第3个二进制位的比较不能区别它们 ,可以将它省略,得到一棵更为简洁的树。



编码: 2: 000010 5: 000101 9: 001001

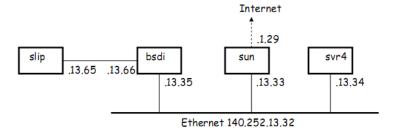
17: 010001 41: 101001 45: 10110163: 111111

压缩Patricia Trie



PATRICIA 应用:路由表查找

Example Net



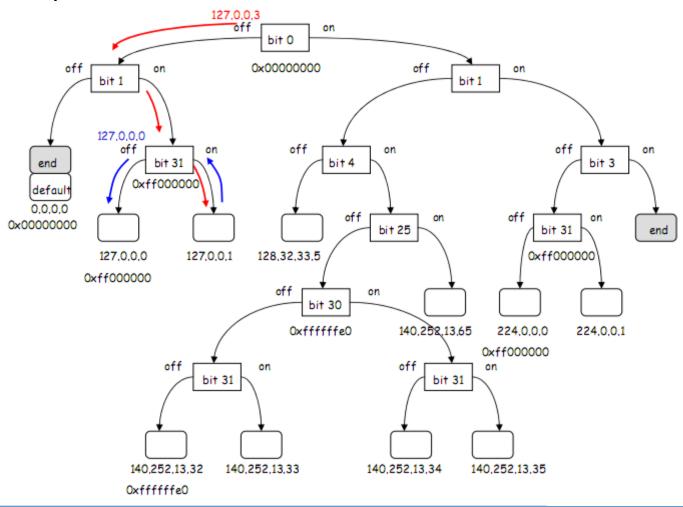
Routing Table

Destination	Gateway	Flags	Ref	Use	Interface
default	140.252.13.33	UGS	0	3	le0
127	127.0.0.1	UGSR	0	2	100
127.0.0.1	127.0.0.1	UH	1	55	100
128.32.33.5	140.252.13.33	UGHS	2	16	le0
140.252.13.32	link#1	UC	0	0	le0
140.252.13.33	8:0:20:3:f6:42	UHL	11	55146	le0
140.252.13.34	0:0:c0:c2:9b:26	UHL	0	3	le0
140.252.13.35	0:0:c0:6f:2d:40	UHL	1	12	100
140.252.13.65	140.252.13.66	UH	0	41	s 10
224	link#1	UC	0	0	le0
224.0.0.1	link#1	UHL	0	5	le0

	32-bit IP address (bits 32-63)	
Bit:	0123 4567 8911 1111 1111 2222 2222 2233 01 2345 6789 0123 4 567 8901	
C	0000 0000 0000 0000 0000 0000 0000 0000 0111 1111 0000 0000	0.0.0.0 127.0.0.0 127.0.0.1 128.32.33.5 140.252.13.32 140.252.13.33 140.252.13.34 140.252.13.35 140.252.13.65 224.0.0.0 224.0.0.1

PATRICIA应用:路由表查找

Example:127.0.0.3



Trie优化(3)--Trie 图

▶ 背景

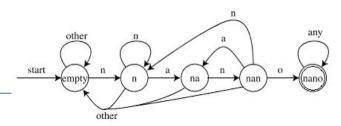
■ 反病毒引擎一般都用字符串作为恶意软件特征,在大规模网络流量中如何判断是否出现恶意软件?

▶ 问题

- □ 给一个很长很长的母串长度为n,然后给m个小模式串。
- □ 求这m个模式串里边有多少个是母串的字串。

➤ Trie图

- □ 一个确定性有限自动机DFA,由Trie树为基础构造
- □ 每点都是一个状态,状态之间的转换用有向边来表示
- □ 用母串作为DFA的输入,在DFA上行走,走到终止节点,就意味着匹配了相应的模式串

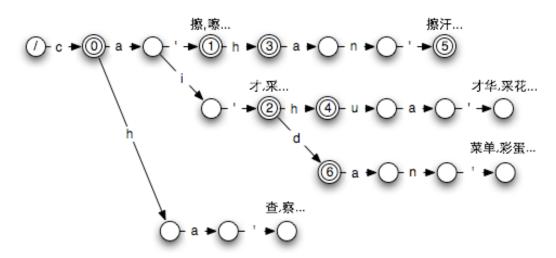


Trie图实例—输入法

举例:

□ 音节表

string[] _spellMusicCode = new string[]{ "a", "ai", "an", "ang", "ao", "ba", "bai", "ban", "ban", "bao", "bei", "ben", "beng", "bi", "bian", "biao", "bie", "bin", "bing", "bo", "bu", "ca", "cai", "can", "cang", "cao", "ce", "ceng", "cha", "chai", "chan", "chang", "che", "chen", "cheng", "chi", "chong", "chou", "chu", "chuai", "chuan", "chuang", "chui", "chun", "chun", "chuan", "chuang", "chui", "chun", "chun", "chuan", "chuang", "chui", "chun", "chuang", "chui", "chun", "chuang", "chui", "chuang", "chui", "chuang", "chui", "chuang", "chui", "chuang", "chui", "chuang", "



本实验的Trie树应用

Ctil	Dec	Hex	Char	Code	Dec	Hex	Char
^@	0	00		NUL	32	20	sp
٩	1	01	©	SOH	33	21	1
°B	2	02	8	SIX	34	22	11
°C	3	03		EIX	35	23	#
٩D	4	04	*	EOI	36	24	\$
۰E	5	05	•	ENQ	37	25	2
۰F	6	06	+	ACK	38	26	&
°G	7	07	•	BEL	39	27	,
°H	8	08	٠	BS	40	28	(
۰I	9	09	0	HI	41	29)
٥J	10	0A	0	LF	42	2A	*
٩K	11	0B	3	VI	43	2B	+
^L	12	0C	Q	FF	44	2C	30
·Μ	13	0D	ſ	CIR	45	2D	-
۰N	14	0E	П	SO.	46	2E	•3
0	15	0F	*	SI	47	2F	
^P	16	10	•	SLE	48	30	0
^Q	17	11	4	CS1	49	31	1
^R	18	12	1	DC2	50	32	2
28	19	13	!!	DC3	51	33	3
٩ī	20	14	q q	DC4	52	34	4
٠U	21	15	δ	NAK	53	35	5
٠V	22	16	= 2	SYN	54	36	6
٠w	23	17	ŧ	EIB	55	37	7
°X	24	18	Ť	CAN	56	38	8
٩Y	25	19	1	EM	57	39	9
۰z	26	1A	->	SIB	58	3A	2
]^	27	1B	+	ESC	59	3B	;
7	28	1C	2	FS	60	3C	1
^]	29	1D	#	GS	61	3D	=
00	30	1E	•	RS	62	3E	>
۰_	31	1F		US	63	3F	?

Dec	Hex	Char	П	Dec	Hex	Char
64	40	6	1	96	60	
65	41	A		97	61	a
66	42	В		98	62	Ъ
67	43	C		99	63	C
68	44	D		100	64	d
69	45	E		101	65	e
70	46	F		102	66	f
71	47	G		103	67	g
72	48	H		104	68	h
73	49	I		105	69	i
74	4 A	J		106	6A	j
75	4B	K		107	6B	k
76	4C	L		108	6C	1
77	4D	M		109	6D	m
78	4 E	N		110	6E	n
79	4 F	0		111	6F	O
30	50	P		112	70	P
31	51	Q		113	71	q
32	52	R		114	72	r
83	53	S		115	73	\$
84	54	T		116	74	t
85	55	U		117	75	u
36	56	V		118	76	v
37	57	W		119	77	w
88	58	X		120	78	×
89	59	Y		121	79	y
90	5A	Z		122	7A	Z
91 (ΔBV	J.E		123	7B	-
92	3C	m	c	124	TC	1
93	5D	1	-	125	7D. 0	om
94	5E	<u>^</u> d		126	李	48
95	5F	1-	1	127	TF	24

Email格式:

- 开头为数字或字母
- 中间可含"_"、"-"、""3类字符
- 不区分大小写
- local part 为"@"前面的部分, 最多64个字符
- domain part 为"@"后面的 部分,最多255个字符
- 总长度不超过 255+64+"@"=320字符

参考实现(Trie)

- 著名实现代码
- <u>Linux Kernel Implementation</u>, used for the page cache, among other things.
- GNU C++ Standard library has a trie implementation
- <u>Java implementation of Concurrent Radix Tree</u>, by Niall Gallagher
- C# implementation of a Radix Tree
- Practical Algorithm Template Library, a C++ library on PATRICIA tries (VC++ >=2003, GCC G++ 3.x), by Roman S. Klyujkov
- <u>Patricia Trie C++ template class implementation</u>, by Radu Gruian
- <u>Haskell standard library implementation</u> "based on big-endian patricia trees".
 Web-browsable source code.
- <u>Patricia Trie implementation in Java</u>, by Roger Kapsi and Sam Berlin
- <u>Crit-bit trees</u> forked from C code by Daniel J. Bernstein
- Patricia Trie implementation in C, in <u>libcprops</u>
- <u>Patricia Trees : efficient sets and maps over integers in OCaml</u>, by Jean-Christophe Filliâtre
- Radix DB (Patricia trie) implementation in C, by G. B. Versiani



Hash思想回顾

- > 思想
 - □ 把任意长度的输入通过Hash算法,变换成固定长度的输出
 - 把一些不同长度的信息转化成杂乱的128/256位的编码
 - 代表算法: MD5, SHA
 - □ 特点
 - 一般是压缩映射,不可逆映射
 - 把一个大范围映射到一个小范围
- > 使用领域
 - □ 信息安全基础设施实现加密、消息摘要等
 - □ 数字指纹:文档去重、内容检索
 - □ 内容安全:数据防泄漏

数据结构之Hash Table

- > 思想
 - □ 利用线性表存储集合元素
 - □ 利用Hash函数计算元素对应的地址entry,并在对应的区域存储 该元素
- ▶ 问题
 - □ 解决冲突
 - □ 加大表空间,
 - 查找复杂(払
 - □ 准确率100%

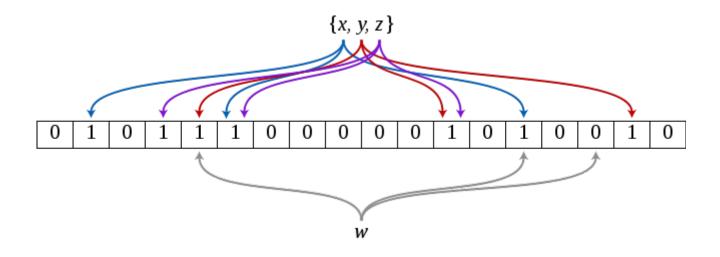
如何改进?

Bloom Filter概念

- ▶ 背景
 - □ 1970年Burton Bloom论文《Space/time trade-offs in hash coding with errors》中提出
 - 吴军《数学之美》中称之为布隆过滤器
- > 概念
 - □ 一个很长的二进制向量和一系列随机映射函数
 - □ 用于检索一个元素是否在一个集合中
 - □ 准确率换空间思想延伸
- ▶ 优点
 - □ 空间效率和查询时间都远超过一般的算法
- > 缺点
 - □ 有一定的误识别率和删除困难

Bloom Filter构建

- > 定义
 - □ 将n个元素集合 $S=\{x_1, x_2,...,x_n\}$
 - □ 一个包含m位的二进制位数组存储
 - □ K个相互独立的哈希函数映射到{1,...,m}的范围
 - □ S集合中的每个元素用k个hash函数映射到, {1,...,m}范围内, 将相应的位置为1



Bloom Filter原理

初始化时m 位数组置零

B 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

对集合中的每个元素 x_i 分别进行k 次hash,If $H_i(x_i) = a$, set B[a] = 1.

B 0 1 0 0 1 0 1 0 1 1 1 1 0 1 0

要检测y是否在集合S,测试所有 $k \land B[H_i(y)]$ 是否都为1.

可能出现false positive: 即所有k个值都是1,但y不在集合S中

n items m = cn bits k hash functions

错误率估计

- ▶ 初始位向量为0
- \blacktriangleright 插入一个元素后,被k个哈希函数(完全独立)映射到位向量后,某一位还是为0的概率是: $\left(1-\frac{1}{m}\right)^k$
- 集合S={x₁, x₂,...,x_n}的所有元素都插入到位向量后,某一位还是为0的概率就是 $p' = \left(1 \frac{1}{m}\right)^{tn} \approx e^{-tn/m}$. $\lim_{x \to \infty} \left(1 \frac{1}{x}\right)^{-x} = e$
- 令 ρ 为 位 向 量 中 0 的 比 例 , 则 ρ 的 数 学 期 望 E(ρ) = p' 。 令 $p = e^{-kn/m}$ 在 ρ 已 知 时 要 求 的 错 误 率 (false positive rate) 为 :

$$(1-\rho)^k \approx (1-p')^k \approx (1-p)^k$$
. $f = (1-e^{-kn/m})^k = (1-p)^k$.

 $(1-\rho)$ 为位数组中1的比例, $(1-\rho)^k$ 就表示k次哈希都刚好选中1的区域

M. Mitzenmacher已经证明位向量中0的比例非常集中地分布在它的数学期望值的附近

最优的哈希函数个数k

▶ 问题

- □ Bloom Filter要用多个哈希函数将集合映射到位向量中,应该 选择几个哈希函数才能使元素查询时的错误率降到最低?
- $f = (1 e^{-kn/m})^k = (1 p)^k$.
- □ 令g = k ln(1 e^{-kn/m}), 让g取到最小, f自然也取到最小
- □ 将g写成

$$g = -\frac{m}{n}\ln(p)\ln(1-p),$$

> 结果

- □ 很容易看出当p = 1/2,也就是 $k = \ln 2 \cdot (m/n)$ 时,g取得最小值
- 在这种情况下,最小错误率f等于(1/2)^k≈ (0.6185)^{m/n}
- □ 另外,注意到p是位数组中某一位仍是0的概率,所以p=1/2 对应着位数组中0和1各一半
- 换句话说,要想保持错误率低,最好让位数组有一半还空着

位向量大小m

问题

- 在不超过一定错误率的情况下, Bloom Filter至少需要多少位 才能表示全集中任意n个元素的集合。假设全集中共有u个元素 ,允许的最大错误率为 ϵ ,求位数组的位数m。
- 一个确定的位向量可以表示的集合个数:
- m位的位数组可以表示的集合个数: $2^m \binom{n+\epsilon(u-n)}{n}$ 全集中n各元素的集合个数: $\binom{u}{n}$
- 要让m位的位数组能够表示所有n个元素的集合,必须有

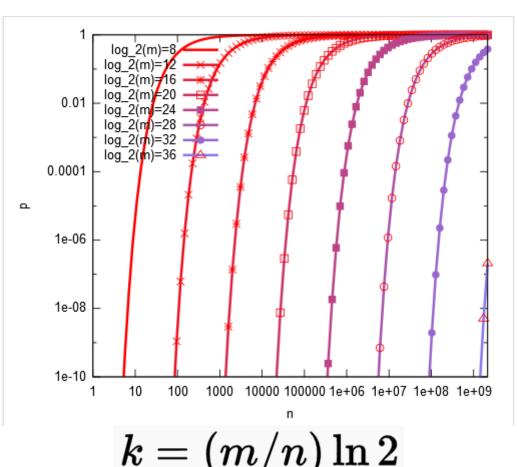
$$2^m \binom{n+\epsilon(u-n)}{n} \ge \binom{u}{n}$$

- 即: $m \ge \log_2 \frac{\binom{u}{n}}{\binom{n+\epsilon(u-n)}{1}} \approx \log_2 \frac{\binom{u}{n}}{\binom{\epsilon u}{n}} \ge \log_2 \epsilon^{-n} = n \log_2 (1/\epsilon).$
- 上页曾算出当 $k = \ln 2 \cdot (m/n)$ 时错误率f最小,这时 $f = (1/2)^k =$ (1/2)^{mln2/n}。现在令f≤e,可以推出,约大了1.44倍

$$m \ge n \frac{\log_2(1/\epsilon)}{\ln 2} = n \log_2 e \cdot \log_2(1/\epsilon).$$

函数分析

False positive rate随Hash 函数个数和存储空间增 加成指数级降低



$$k=(m/n)\ln 2$$

测试与参数选择

- ▶ 进行3组实验,每组取5个N
 - □ 取FP1 =0.01%, N=[50W, 100W, 300W, 500W, 1000W]

 - □ 取FP3 =0.00001%,N=[50W, 100W, 300W, 500W,1000W]

	(Vector size)m			内存			(Hash num) k			Х		
N	FP1	FP2	FP3	FP1	FP2	FP3	FP1	FP2	FP3	FP1	FP2	FP3
50W	958W	1198 W	1677W	1M	1M	1M	13	17	23	37091	3691	38
100W	1917 W	2396 W	3355W	2M	2M	3M	13	17	23	36958	3770	47
300W	5751 W	7188 W	10064 W	6M	8M	11M	13	17	23	36585	3689	38
500W	9585 W	11981 W	16773 W	11 M	14M	19M	13	17	23	36569	3701	45
1000W	19170W	23962W	33547 W	22 M	28M	39M	13	17	23	36533	3552	41

Hash算法选择

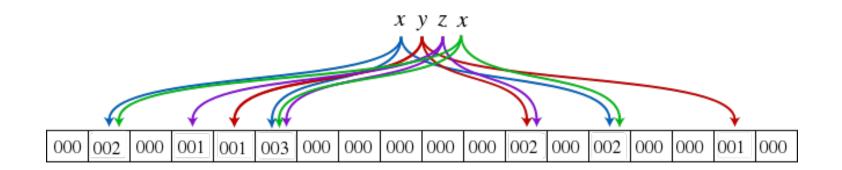
- ► K个hash越相互独立效果越好
- ▶ 常用的HASH算法
 - http://www.partow.net/programming/hashfunctions/
 - unsigned int RSHash (char* str, unsigned int len);
 - unsigned int JSHash (char* str, unsigned int len);
 - unsigned int PJWHash (char* str, unsigned int len);
 - unsigned int ELFHash (char* str, unsigned int len);
 - unsigned int BKDRHash(char* str, unsigned int len);
 - unsigned int SDBMHash(char* str, unsigned int len);
 - unsigned int DJBHash (char* str, unsigned int len);
 - unsigned int DEKHash (char* str, unsigned int len);
 - unsigned int BPHash (char* str, unsigned int len);
 - unsigned int FNVHash (char* str, unsigned int len);
 - unsigned int APHash (char* str, unsigned int len);

Hash算法不够

- ➤ Hash算法特性
 - □ 抗冲突性(collision-resistant),即在统计上无法产生2个散列值相同的预映射。
 - □ 映射分布均匀性和差分分布均匀性
 - ▶ 散列结果中,为0的bit和为1的bit,其总数应该大致相等;
- ➤ 超级HASH算法
 - The requirement of designing k different independent hash functions can be prohibitive for large k. For a good hash function with a wide output, there should be little if any correlation between different bit-fields of such a hash, so this type of hash can be used to generate multiple "different" hash functions by slicing its output into multiple bit fields.
 - MurmurHash
 - □ 官网: https://sites.google.com/site/murmurhash/

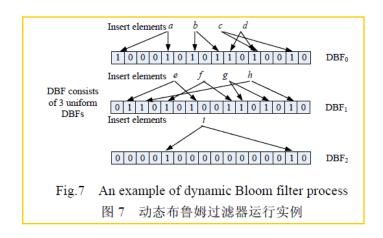
扩展: 计数型Bloom Filter

- ► 标准BF
 - □ 没法删除元素
 - □ m长数组,每unit用1位表示,只能表示[0,1]
- ▶ 如果每个unit用多位表示,如3位,则可以表示[0,1,...,7]



扩展: 压缩Bloom Filter

▶ 由于标准Bloom Filter在f最小时,任意位为0的概率为 1/2, 所以,为了更好的在网络上传输Bloom Filter,可以对Bloom Filter进行压缩,能得到约69.3%





实验输入和输出

- ▶ 输入数据文件
 - □ 总共2000万个字符串,每个字符串一行

```
113247 timothy154473y4002@yahoo.com
113248 timothy154473y5@hotpop.com
113249 timothy154473y741@yahoo.com
113250 timothy15447fds@imailbox.com
113251 timothy15447fds@imailbox.com
113252 timothy15447fds@flashmail.com
113253 timothy15447sdf4796@yahoo.com
113254 timothy15447sdf@imailbox.com
113255 timothy15447sdfg@imailbox.com
113256 timothy15447sdfg@imailbox.com
113257 timothy15447sdfgg5017@yahoo.com
113258 timothy15447sq3413@yahoo.com
113258 timothy15447y365y@hotpop.com
```

- □ 待检测数据1000个字符串,每个字符串一行
- ▶ 输出要求
 - □ 每行输出一个1000行结果的文件,每行一个yes/no,表示 1000个email是否在2000万字符串中

实验说明

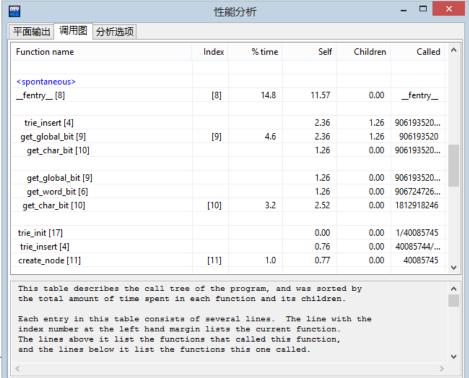
- > 代码文件
 - □ 代码文件
 - main.c bloom.c bloom.h trie.c trie.h hash.c hash.h
 - Makefile编译规则文件
 - Readme说明程序的编译和使用方法
 - □ 命令行格式
 - ./strsearch strlist.txt checklist.dat checkedresult.dat
 - □ 输入文件
 - strlist.dat, 2000万个源字符串
 - checklist.dat, 1000个待检测的字符串
 - □ 输出文件(要求程序生成)
 - checkresult.dat
- > 文档
 - □ 实验XXXXX设计文档.doc

编程实验环境

- 代码和数据服务器
 - □ 见邮件
- > 操作系统
 - □ 推荐Ubuntu Linux操作系统(虚拟机或cygwin)
- > 编辑器
 - Source Insight或Dev C++
- > 编译和调试
 - □ GCC/G++
 - □ gdb
 - Makefile

性能优化工具

- > 性能瓶颈分析
 - □ 可以采用gcc编译工具包gprofile工具来动态分析程序的性能瓶颈
 - □ 可以看到程序的每个函数的调用次数以及占用CPU的时间比例
 - □ Dev-CPP集成了Profile工具,可以直接用可视化界面查看



几个性能优化的方法

- > 函数内联
 - □ 把一些大量调用的方法内联,添加inline关键字

```
inline
ull get_word_bit(char * word, ull index){
    return get_char_bit(word[index / BITS_USED], index % BITS_USED);
}
```

- ▶ 减少malloc次数
 - □ 内存预分配
 - □ tips: tcmalloc (多线程) / Jemalloc
 - 比glibc的malloc速度快2倍以上
 - Thread-Caching Malloc by Google

