## 15-853:Algorithms in the Real World

**Indexing and Searching I**
– Introduction
– Inverted Indices

## Outline for next few classes

**Inverted Indices (used by all search engines)**
– Compression
– The lexicon
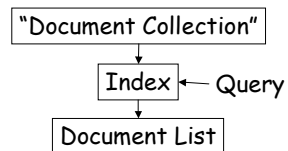– Merging terms (unions and intersections)
**Vector Models**
**Latent Semantic Indexing**
**Link Analysis:**
– PageRank (Google)
– HITS
**Duplicate Removal**

## Basic Model

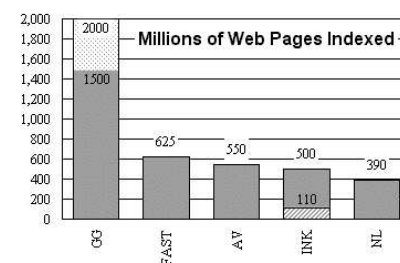"Document Collection"
↓
Index ← Query
↓
Document List

**Applications**:
– Web, mail and dictionary searches
– Law and patent searches
– Information filtering (e.g., NYT articles)
**Goal**: **Speed**, **Space**, **Accuracy**, **Dynamic Updates**

## How big is an Index?



Millions of Web Pages Indexed: GG 2000/1500, FAST 625, AV 550, INK 500/110, NL 390

Dec 2001, self proclaimed sizes (gg = google)
Source: Search Engine Watch

1

## Main Approaches

**%ull text searching**
- e.g. grep, agrep (used by many mailers)

**Inverted Indices**
- good for short >ueries
- used by most search engines

**Signature %iles**
- good for longer >ueries with many terms

**Vector Space Models**
- good for better accuracy
- used in clustering, S?D, @

---

## Queries

**&ypes o' ( ueries on Multiple )terms***
- boolean (and, or, not, andnot)
- proximity (adA within BnO)
- keyword sets
- in relation to other documents

**And +ithin each term**
- prefix matches
- wildcards
- edit distance bounds

---

## Techni>ue used Across Methods

**, ase ' olding**
London -Clondon

**Stemming**
compress = compression = compressed
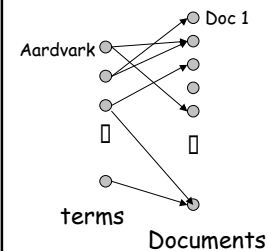(several off-the-shelf English Language stemmers are freely available)

**Stop +ords**
to, the, it, be, or, @
how about "to be or not to be"

**&hesaurus**
fast -Crapid

---

## Documents as Bipartite Graph



Doc 1
Aardvark
@          @
terms
Documents

**, alled an )Inverted %ile Index***

**, an be stored using ad-acency lists, also called**
- posting lists (or files)
- inverted file entry

**. xample si/e o' &R. ,**
- 538D terms
- EF2D documents
- 333,85GD edges

**%or the +eb, multiply by** 01234

2

## Documents as Bipartite Graph



terms

Documents

**Implementation Issues**:

25 **Space ' or posting lists**

these take almost all the space

65 **Access to lexicon**

- btrees, tries, hashing
- prefix and wildcard >ueries

75 **Merging posting list**

- multiple term >ueries

---

## 1. Space for Posting Lists

8**osting lists can be as large as the document data**

- saving space and the time to access the space is critical for performance

9 **e can compress the lists,**

**but, +e need to uncompress on the ' ly**5

**Di' ' erence encoding**:

Lets say the term <u>elephant</u> appears in documents:

:7, 0, 63, 62, 67, ; <, ; ; , ; =>

then the difference code is

:7, 6, 20, 2, 6, 07, 2, 2>

---

## Some Codes

**Gamma code**:

if most significant bit of n is in location k, then

gamma(n) = $0^{k-1}$ nHk..0I

$2 \log(n) - 1$ bits

**Delta code**:

gamma(k)nHk..0I

$2 \log(\log(n))$ J $\log(n) - 1$ bits

**%re?uency coded**:

base on actual probabilities of each distance

---

## Global vs. Local Probabilities

**Global**:

- Count K of occurneces of each distance
- L se Huffman or arithmetic code

**Local**:

generate counts for each list

elephant: H3, 2, 1, 2, 53, 1, 1I

Problem: counts take too much space

Solution: batching

group into buckets by $\lfloor \log(\text{length}) \rfloor$

## Performance

| Global | bitsEedge |
|---|---|
| Binary | 63533 |
| Gamma | <5D7 |
| Delta | <52C |
| Huffman | 05=7 |
| **Local** | |
| Skewed Bernoulli | 056= |
| Batched Huffman | 056; |

**@its per edge based on the &R. , document collection**

**&otal si/e** A 777**M** B 5<< **bytes** A 666**Mbytes**

---

## 2. Accessing the Lexicon

9 **e all kno+ ho+ to store a dictionary, @U&F**
- it is best if lexicon fits in memory---can we avoid storing all characters of all words
- what about prefix or wildcard >ueries?

**Some possible data structures**
- M´ont Coding
- Tries
- Perfect Hashing
- B-trees

---

## M´ont Coding

| 9 ord | ' ront coding |
|---|---|
| E,Ӕzebel | 0,E,Ӕzebel |
| 5,Ӕzer | F,1,r |
| E,Ӕzerit | 5,2,it |
| G,Ӕziah | 3,3,iah |
| G,Ӕziel | F,2,el |
| E,Ӕzliah | 3,F,liah |

**%or large lexicons can save** ; 0G **o' space**

**@ut +hat about random access**H

---

## Prefix and Wildcard Queries

8**re' ix ?ueries**
- Handled by all access methods except hashing
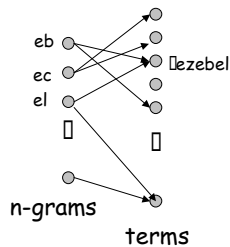
9 **ildcard ?ueries**
- n-gram
- rotated lexicon

4

# n-gram

**, onsider every block o' n characters in a term:**

e.g. 2-gram of Æzebel -C$j,je,ez,ze,eb,el,l$

eb
ec
el
@

n-grams

Æzebel
@

terms

Break wildcard >uery into an n-grams and search.

e.g. `j*el` would

1. search for `$j,el,l$` as if searching for documents
2. find all potential terms
3. filter matches for which the order does not match

---

# Rotated Lexicon

**, onsider every rotation o' a term:**

e.g. Æzebel -CNÆezebel, lNÆezebe, elNÆezeb, belNÆeze

**I o+ store lexicon o' all rotations**

**Given a ?uery ' ind longest contiguous block (+ith rotation) and search ' or it:**

e.g. AŒl -Csearch for elNAin lexicon

**I ote that each lexicon entry corresponds to a single term**

e.g. ebelNÆez can only mean Æzebel

---

# 3. Merging Posting Lists

**Lets say ?ueries are expressions over:**
– and, or, andnot

**Vie+ the list o' documents ' or a term as a set:**

**&hen**

$e_1$ and $e_2$ -C$S_1$ intersect $S_2$

$e_1$ or $e_2$ -C$S_1$ union $S_2$

$e_1$ andnot $e_2$ -C$S_1$ diff $S_2$

**Some notes:**
– the sets ordered in the "posting lists"
– $S_1$ and $S_2$ can differ in size substantially
– might be good to keep intermediate results
– persistence is important

---

# L nion, Intersection, Merging

**Given t+o sets o' length $n$ and $m$ ho+ long does it take ' or intersection, union and set di' ' erenceH**

**Assume elements are taken ' rom a total order (J)**

**Very similar to merging t+o sets A and @, ho+ long does this takeH**

**Lo+er @ound:**
– There are n elements of A and n J m positions in the output they could belong
– choose (n J m, n) possibilities
– assuming comparison based model, the decision tree has that many leaves and depth log of that
– Assuming m Bn this give $\Omega(m \log((n J m)Pn))$

## Merging: Upper bounds

Tarjan shows $O(m \log((n+m)/n))$ upper bounds using 2-3 trees with cross links and parent pointers. Very messy.

We will take different approach, and base on two operations: **split** and **-oin**

© Guy Blelloch, 2002          15-853          Page21

---

## Split and Qoin

**Split(S,v)** : Split S into two sets $S_B = \{s \in S \mid s < v\}$ and $S_C = \{s \in S \mid s > v\}$. Also return a flag which is true if $v \in S$.
- Split({5,9,15,18,22}, 18) → {5,9,15},{22},True

**Koin(S_B, S_C)** : Assuming $\forall k_B \in S_B, k_C$ in $S_C: k_B < k_C$ it returns $S_B \cup S_C$
- Qoin({5,9,11},{17,22}) → {5,9,11,17,22}

**Time for both:**
- $O(\log(\min(|S_B|, |S_C|)))$, can be shown
- $O(\log |S_B|)$, will suffice for us (shown later)

© Guy Blelloch, 2002          15-853          Page22

---

## Union with Split and Qoin

**Union(S$_1$, S$_2$)** =
   if isempty(S$_1$) **then return** S$_2$
   **else**
      (S$_{2B}$, S$_{2C}$, fl) = Split(S$_2$, first(S$_1$))
      **return** Qoin(S$_{2B}$, Union(S$_{2C}$, S$_1$))

| A | a1 | a2 | a3 |  | aF |  | a5 |
|---|----|----|----|--|----|--|----|

| B | b1 | b2 | b3 | bF | b5 |
|---|----|----|----|----|----|

| Out | b1 | a1 | b2 | a2 | b3 | a3 | bF | aF |
|-----|----|----|----|----|----|----|----|----|

© Guy Blelloch, 2002          15-853          Page23

---

## Runtime of Union

| Out | o1 | o2 | o3 | oF | o5 | oG | oE | o8 |
|-----|----|----|----|----|----|----|----|----|

$T_{union} = O(\sum_i \log |o_i| + \sum_i \log |o_i|)$
               **Splits        Koins**

Since the logarithm function is concave, this is maximized when blocks are as close as possible to equal size, therefore

$T_{union} = O(\sum_{i=1}^{m} \log \lceil n/m + 1 \rceil)$
             $= O(m \log ((n+m)/m))$

© Guy Blelloch, 2002          15-853          Page24

6

## Intersection with Split and Qoin

**Intersect(S$_2$, S$_2$)** A
   i' isempty(S$_1$) **then return** $\varnothing$
   **else**
      (S$_{2B}$, S$_{2G}$ flag) = Split(S$_2$, first(S$_1$))
      i' flag **then**
        **return** Qoin(first(S$_1$), Intersect(S$_{2G}$ S$_1$))
      **else**
        **return** Intersect(S$_{2G}$ S$_1$)

## Efficient Split and Qoin

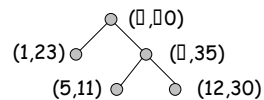**Recall that +e +ant:** T = $O(\log \frac{S S_B}{S})$

How do we implement this efficiently?

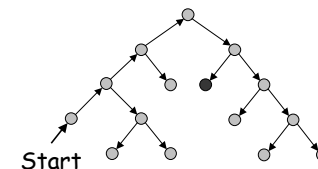## Treaps

. **very key is given a )random* priority**5
  – keys are stored in-order
  – priorities are stored in heap-order
e.g. (key,priority) : (1,23), (F,F0), (5,11), (U,35), (12,30)



If the priorities are uni>ue, the tree is uni>ue.

## Left Spinal Treap



Start

**&ime to split** A **length** ' **rom Start to split location**
9 **e +ill sho+ that this is** L **(log L) in the expected case, +here L is the path length bet+een Start and the split location**
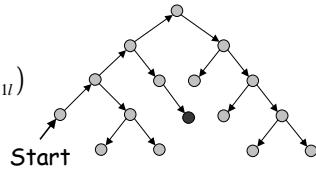**&ime to** Koin **is the same**

## Analysis

$P_i = $ lenght of path from Start to $i$      $p_i = Ex[P_i]$

$$A_{ij} = \begin{cases} 1 & x_i \text{ ancestor of } x_j \\ 0 & \text{otherwise} \end{cases}$$      $a_{ij} = Ex[A_{ij}]$

$$C_{ilm} = \begin{cases} 1 & x_i \text{ common ancestor of } x_l \text{ and } x_m \\ 0 & \text{otherwise} \end{cases}$$    $c_{ilm} = Ex[C_{ilm}]$

$$P_l = \sum_{i=1}^{l} A_{i1} + \sum_{i=1}^{n} \left(A_{il} - C_{i1l}\right)$$

Start

## Analysis Continued

$$Ex[P_l] = p_l = \sum_{i=1}^{l} a_{i1} + \sum_{i=1}^{n} \left(a_{il} - c_{i1l}\right)$$

**Lemma**: $a_{ij} = \dfrac{1}{|i-j|+1}$

8**roo**':

1. i is an ancestor of A iff i has a greater priority than all elements between i and A inclusive.
2. there are S-ASJ1 such elements each with e>ual probability of having the highest priority.

## Analysis Continued

$$\sum_{i=1}^{l} a_{i1} = \sum_{i=1}^{l} \frac{1}{|i-1|+1} = \sum_{i=1}^{l} \frac{1}{i}$$
$$< 1 + \ln l \quad \text{(harmonic number } H_l\text{)}$$

Can similarly show that:

$$\sum_{i=1}^{n} \left(a_{il} - c_{i1l}\right) = O(\log l)$$

Therefore the epected path length and runtime for split and Join is O(log l).

Similar techni>ue can be used for other properties of Treaps.

## And back to Inverted Indices

9 **e sho+ed ho+ to take Unions and Intersections, but** &**reaps are not very space e''icient**5

**Idea: i' priorities are in the range** :3552**) then any node +ith priority** J 2 1 **α is stored compressed**5

**α represents 'raction o' uncompressed nodes**5

8