# CSE 454

Index Compression
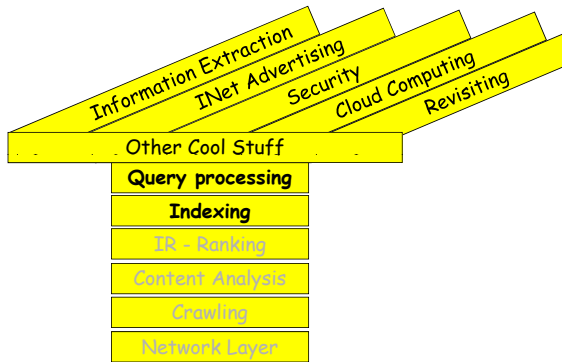Alta Vista
PageRank

---

# Administrivia

- **No class Tues 10/26**
  - Instead go to today's colloquium

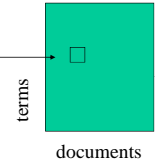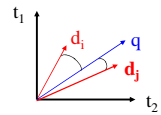  - Group Meetings

- **Never-Ending Language Learning**
  - Today 3:30pm EEB 105

---

# Class Overview

Information Extraction
INet Advertising
Security
Cloud Computing
Revisiting

Other Cool Stuff

**Query processing**

**Indexing**

IR - Ranking

Content Analysis

Crawling

Network Layer

---

# Review

- **Vector Space Representation**
  - Dot Product as Similarity Metric

- **TF-IDF for Computing Weights**
  - $w_{ij} = f(i,j) * log(N/n_i)$
  - Where $q = \ldots word_i \ldots$
  - $N = |docs|$     $n_i = |docs\ with\ word_i|$

- **But How Process Efficiently?**

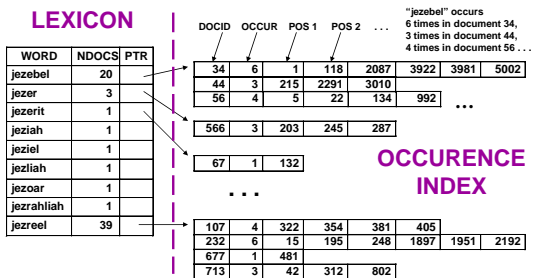$t_1$    $d_i$    q    $d_j$    $t_2$

terms

documents

---

# Retrieval

Document-term matrix

|       | $t_1$ | $t_2$ | $\ldots$ | $t_j$ | $\ldots$ | $t_m$ | nf |
|-------|-------|-------|------|-------|------|-------|------|
| $d_1$ | $w_{11}$ | $w_{12}$ | $\ldots$ | $w_{1j}$ | $\ldots$ | $w_{1m}$ | $1/|d_1|$ |
| $d_2$ | $w_{21}$ | $w_{22}$ | $\ldots$ | $w_{2j}$ | $\ldots$ | $w_{2m}$ | $1/|d_2|$ |
| $d_i$ | $w_{i1}$ | $w_{i2}$ | $\ldots$ | $w_{ij}$ | $\ldots$ | $w_{im}$ | $1/|d_i|$ |
| $d_n$ | $w_{n1}$ | $w_{n2}$ | $\ldots$ | $w_{nj}$ | $\ldots$ | $w_{nm}$ | $1/|d_n|$ |

$w_{ij}$ is the weight of term $t_j$ in document $d_i$
Most $w_{ij}$'s will be zero.

---

# Inverted Files for Multiple Documents

**LEXICON**

"jezebel" occurs
6 times in document 34,
3 times in document 44,
4 times in document 56 . . .

| WORD | NDOCS | PTR |
|------|-------|-----|
| jezebel | 20 | |
| jezer | 3 | |
| jezerit | 1 | |
| jeziah | 1 | |
| jeziel | 1 | |
| jezliah | 1 | |
| jezoar | 1 | |
| jezrahliah | 1 | |
| jezreel | 39 | |

| DOCID | OCCUR | POS 1 | POS 2 | $\ldots$ | | | |
|-------|-------|-------|-------|------|---|---|---|
| 34 | 6 | 1 | 118 | 2087 | 3922 | 3981 | 5002 |
| 44 | 3 | 215 | 2291 | 3010 | | | |
| 56 | 4 | 5 | 22 | 134 | 992 | $\ldots$ | |
| 566 | 3 | 203 | 245 | 287 | | | |
| 67 | 1 | 132 | | | | | |

**OCCURENCE INDEX**

. . .

| 107 | 4 | 322 | 354 | 381 | 405 | | |
| 232 | 6 | 15 | 195 | 248 | 1897 | 1951 | 2192 |
| 677 | 1 | 481 | | | | | |
| 713 | 3 | 42 | 312 | 802 | | | |

## Many Variations Possible

- **Address space (flat, hierarchical)**
  - Alta Vista uses flat approach
- **Record term-position information**
- **Precalculate TF-IDF info**
- **Stored header, font & tag info**
- **Compression strategies**

7

---

## Compression

- **What Should We Compress?**
  - Repository
  - Lexicon
  - Inv Index
- **What properties do we want?**
  - Compression ratio
  - Compression speed
  - Decompression speed
  - Memory requirements
  - Pattern matching on compressed text
  - Random access

8

---

## Inverted File Compression

Each inverted list has the form $< f_t ; d_1 , d_2 , d_3 , ... , d_{f_t} >$

A naïve representation results in a storage overhead of $(f + n) * \lceil \log N \rceil$

This can also be stored as $< f_t ; d_1 , d_2 - d_1 , ... , d_{f_t} - d_{f_t - 1} >$

Each difference is called a d-gap. Since $\sum (d - gaps) \leq N$,

each pointer requires fewer than $\lceil \log N \rceil$ bits.

Trick is encoding …. since worst case ….

➡ *Assume d-gap representation for the rest of the talk, unless stated otherwise*

Slides adapted from Tapas Kanungo and David Mount, Univ Maryland

9

---

## Text Compression

Two classes of text compression methods
- Symbolwise (or statistical) methods
  - **Estimate probabilities of symbols - modeling step**
  - **Code one symbol at a time - coding step**
  - **Use shorter code for the most likely symbol**
  - **Usually based on either arithmetic or Huffman coding**
- Dictionary methods
  - **Replace fragments of text with a single code word**
  - **Typically an index to an entry in the dictionary.**
    - **eg: Ziv-Lempel coding: replaces strings of characters with a pointer to a previous occurrence of the string.**
  - **No probability estimates needed**

➡ *Symbolwise methods are more suited for coding d-gaps*

10

---

## Classifying d-gap Compression Methods:

- **Global: each list compressed using same model**
  - **non-parameterized**: probability distribution for d-gap sizes is predetermined.
  - **parameterized**: probability distribution is adjusted according to certain parameters of the collection.
- **Local: model is adjusted according to some parameter, like the frequency of the term**

- **By definition, local methods are parameterized.**

11

---

## Conclusion

- **Local methods best**

- **Parameterized global models ~ non-parameterized**
  - Pointers not scattered randomly in file
- **In practice, best index compression algorithm is:**
  - Local Bernoulli method (using Golomb coding)
- **Compressed inverted indices usually faster+smaller than**
  - Signature files
  - Bitmaps

Local < Parameterized Global < Non-parameterized Global

Not by much

12

## CSE 454 - Case Studies

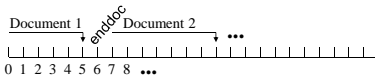### Design of Alta Vista

**Based on a talk by Mike Burrows**

## AltaVista: Inverted Files

- **Map each word to list of locations where it occurs**
- **Words = null-terminated byte strings**
- **Locations = 64 bit unsigned ints**
  - Layer above gives interpretation for location
    - URL
    - Index into text specifying word number

- **Slides adapted from talk by Mike Burrows**

## Documents

- **A document is a region of location space**
  - Contiguous
  - No overlap
  - Densely allocated (first doc is location 1)
- **All document structure encoded with words**
  - enddoc at last location of document
  - begintitle, endtitle mark document title

```
Document 1   enddoc  Document 2   •••
 |_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|_|
 0 1 2 3 4 5 6 7 8  •••
```

## Format of Inverted Files

- **Words ordered lexicographically**
- **Each word followed by list of locations**
- **Common word prefixes are compressed**
- **Locations encoded as deltas**
  - Stored in as few bytes as possible
  - 2 bytes is common
  - Sneaky assembly code for operations on inverted files
    - Pack deltas into aligned 64 bit word
    - First byte contains continuation bits
    - Table lookup on byte => no branch instructs, no mispredicts
    - 35 parallelized instructions/ 64 bit word = 10 cycles/word
- **Index ~ 10% of text size**

## Index Stream Readers (ISRs)

- **Interface for**
  - Reading result of query
  - Return ascending sequence of locations
  - Implemented using lazy evaluation
- **Methods**
  - loc(ISR)              return current location
  - next(ISR)             advance to next location
  - seek(ISR, X)          advance to next loc after X    ⬅ !
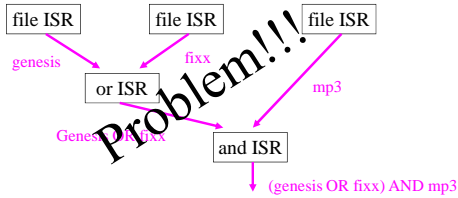  - prev(ISR)             return previous location

## Processing Simple Queries

- **User searches for "mp3"**

- **Open ISR on "mp3"**
  - Uses hash table to avoid scanning entire file
- **Next(), next(), next()**
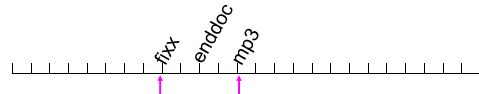  - returns locations containing the word

## Combining ISRs

- **And**     **Compare locs on two streams**
- **Or**       **Merges two or more ISRs**
- **Not**      **Returns locations not in ISR (lazily)**

```
file ISR        file ISR        file ISR
```
genesis          fixx
```
or ISR                          mp3
```
**Problem!!!**

Genesis OR fixx
```
                    and ISR
```
                    (genesis OR fixx) AND mp3

## What About File Boundaries?

fixx  enddoc  mp3

## ISR Constraint Solver

- **Inputs:**
  – Set of ISRs:  A, B, ...
  – Set of Constraints
- **Constraint Types**
  – $loc(A) \leq loc(B) + K$
  – $prev(A) \leq loc(B) + K$
  – $loc(A) \leq prev(B) + K$
  – $prev(A) \leq prev(B) + K$
- **For example: phrase "a b"**
  – $loc(A) \leq loc(B)$,   $loc(B) \leq loc(A) + 1$

a   a   b a     a b      a      b

## Two words on one page

- **Let E be ISR for word** enddoc
- **Constraints for conjunction a AND b**
  – $prev(E) \leq loc(A)$
  – $loc(A) \leq loc(E)$
  – $prev(E) \leq loc(B)$
  – $loc(B) \leq loc(E)$

*What if prev(E) Undefined?*

b   a   e b      a b        e        b

prev(E)              loc(B)     loc(E)

loc(A)

## Advanced Search

- **Field query:   a in Title of page**
- **Let BT, ET be ISRP of words** begintitle, endtitle
- **Constraints:**
  – $loc(BT) \leq loc(A)$
  – $loc(A) \leq loc(ET)$
  – $prev(ET) \leq loc(BT)$

et   a       bt   a   et        a         bt   et

prev(ET)

loc(A)   loc(ET)

loc(BT)

## Implementing the Solver

**Constraint Types**
  – $loc(A) \leq loc(B) + K$
  – $prev(A) \leq loc(B) + K$
  – $loc(A) \leq prev(B) + K$
  – $prev(A) \leq prev(B) + K$

4

## Slide 25

# *Remember:* Index Stream Readers

- **Methods**
  - loc(ISR)        return current location
  - next(ISR)       advance to next location
  - seek(ISR, X)    advance to next loc after X
  - prev(ISR)       return previous location

## Slide 26

# Solver Algorithm

| | |
|---|---|
| loc(ISR) | return cur loc |
| next(ISR) | adv to nxt loc |
| | return it |
| seek(ISR, X) | adv to nxt loc |
| | return it |
| prev(ISR) | return pre loc |

while (unsatisfied_constraints)
    satisfy_constraint(choose_unsat_constraint())

- **To satisfy:    $loc(A) \leq loc(B) + K$**
  - Execute: seek(B, loc(A) - K)

## Slide 27

# Solver Algorithm

| | |
|---|---|
| loc(ISR) | return cur loc |
| next(ISR) | adv to nxt loc |
| | return it |
| seek(ISR, X) | adv to nxt loc |
| | return it |
| prev(ISR) | return pre loc |

while (unsatisfied_constraints)
    satisfy_constraint(choose_unsat_constraint())

- **To satisfy:    $loc(A) \leq loc(B) + K$**
  - Execute: seek(B, loc(A) - K)
- **To satisfy:    $prev(A) \leq loc(B) + K$**
  - Execute: seek(B, prev(A) - K)

## Slide 28

# Solver Algorithm

| | |
|---|---|
| loc(ISR) | return cur loc |
| next(ISR) | adv to nxt loc |
| | return it |
| seek(ISR, X) | adv to nxt loc |
| | return it |
| prev(ISR) | return pre loc |

while (unsatisfied_constraints)
    satisfy_constraint(choose_unsat_constraint())

- **To satisfy:    $loc(A) \leq loc(B) + K$**
  - Execute: seek(B, loc(A) - K)
- **To satisfy:    $prev(A) \leq loc(B) + K$**
  - Execute: seek(B, prev(A) - K)
- **To satisfy:    $loc(A) \leq prev(B) + K$**
  - Execute: seek(B, loc(A) - K),
  -         next(B)

## Slide 29

# Solver Algorithm

| | |
|---|---|
| loc(ISR) | return cur loc |
| next(ISR) | adv to nxt loc |
| | return it |
| seek(ISR, X) | adv to nxt loc |
| | return it |
| prev(ISR) | return pre loc |

while (unsatisfied_constraints)
    satisfy_constraint(choose_unsat_constraint())

- **To satisfy:    $loc(A) \leq loc(B) + K$**
  - Execute: seek(B, loc(A) - K)
- **To satisfy:    $prev(A) \leq loc(B) + K$**
  - Execute: seek(B, prev(A) - K)
- **To satisfy:    $loc(A) \leq prev(B) + K$**
  - Execute: seek(B, loc(A) - K),
  -         next(B)
- **To satisfy:    $prev(A) \leq prev(B) + K$**
  - Execute: seek(B, prev(A) - K)
  -         next(B)

## Slide 30

# Solver Algorithm

Heuristic: Which choice advances a stream the furthest?

while (unsatisfied_constraints)
    satisfy_constraint(choose_unsat_constraint())

- **To satisfy:    $loc(A) \leq loc(B) + K$**
  - Execute: seek(B, loc(A) - K)
- **To satisfy:    $prev(A) \leq loc(B) + K$**
  - Execute: seek(B, prev(A) - K)
- **To satisfy:    $loc(A) \leq prev(B) + K$**
  - Execute: seek(B, loc(A) - K),
  -         next(B)
- **To satisfy:    $prev(A) \leq prev(B) + K$**
  - Execute: seek(B, prev(A) - K)
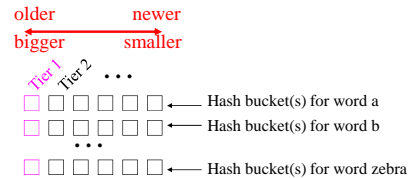  -         next(B)

## Update

- **Can't insert in the middle of an inverted file**
- **Must rewrite the entire file**
  - Naïve approach: need space for two copies
  - Slow since file is huge
- **Split data along two dimensions**
  - **Buckets** solve disk space problem
  - **Tiers** alleviate small update problem

## Buckets & Tiers
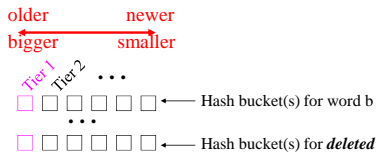
- **Each word is hashed to a bucket**
- **Add new documents by adding a new tier**
  - Periodically merge tiers, bucket by bucket

older → newer
bigger → smaller

Tier 1  Tier 2  . . .

☐ ☐ ☐ ☐ ☐ ← Hash bucket(s) for word a
☐ ☐ ☐ ☐ ☐ ← Hash bucket(s) for word b
. . .
☐ ☐ ☐ ☐ ☐ ← Hash bucket(s) for word zebra

## What if Word Removed from Doc?

- **Delete documents by adding** deleted **word**

deleted ↓
☐☐☐☐☐☐ ▬▬▬▬ ☐☐☐☐☐☐☐☐☐☐☐
↑
a

- **Expunge deletions when merging tier 1**

older → newer
bigger → smaller

Tier 1  Tier 2  . . .

☐ ☐ ☐ ☐ ☐ ← Hash bucket(s) for word b
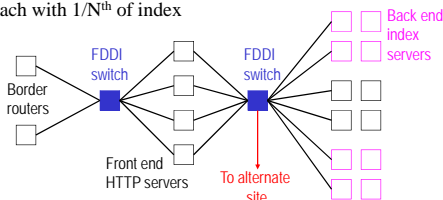. . .
☐ ☐ ☐ ☐ ☐ ← Hash bucket(s) for *deleted*

## Scaling

- **How handle huge traffic?**
  - AltaVista Search ranked #16
  - 10,674,000 unique visitors (Dec'99)
- **Scale across N hosts**
  1. Ubiquitous index. Query one host
  2. Split N ways. Query all, merge results
  3. Ubiquitous index. Host handles subrange of locations. Query all, merge results
  4. Hybrids

## AltaVista Structure

- **Front ends**
  - Alpha workstations
- **Back ends**
  - 4-10 CPU Alpha servers
    - 8GB RAM, 150GB disk
  - Organized in groups of 4-10 machines
    - Each with 1/$N^{th}$ of index

Border routers

FDDI switch

Front end HTTP servers

FDDI switch

To alternate site

Back end index servers