

CSE 454

Search Engine Query Processing Alta Vista Case Study

Based on a talk by Mike Burrows

2/5/2013 2:21 PM Copyright © 2000-2009 D.S.Weld 1

Review

- Vector Space Representation**
 - Dot Product as Similarity Metric
- TF-IDF for Computing Weights**
 - $w_{ij} = f(i,j) * \log(N/n_i)$
 - Where $q = \dots \text{word}_i \dots$
 - $N = |\text{docs}|$ $n_i = |\text{docs with word}_i|$
- But How Process Efficiently?**

Copyright © 2000-2009 D.S.Weld 2

Inverted Files for Multiple Documents

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

LEXICON

DOCID OCCUR POS 1 POS 2 ...

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

34	6	1	118	2087	3922	3981	5002
44	3	215	2291	3010			
56	4	5	22	134	992		
...							
566	3	203	245	287			
67	1	132					
...							
107	4	322	354	381	405		
232	6	15	195	248	1897	1951	2192
677	1	481					
713	3	42	312	802			

OCCURENCE INDEX

Copyright © 2000-2009 D.S.Weld 4

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**

How Process Complex Queries?

- Phrases, term in anchor text, etc**

Copyright © 2000-2009 D.S.Weld 5

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
- Preview:**
 - **Index Stream Readers:** four primitives: loc, seek, ...
 - **Constraint Solvers:** four constraints: $loc_A \leq loc_B + K$, ...

2/5/2013 2:21 PM 6

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

2/5/2013 2:21 PM 7

Format of Inverted Files

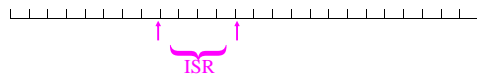
- Words ordered lexicographically
- Each word followed by list of locations
- Common word prefixes are compressed
- Locations encoded as deltas
 - Stored in as few bytes as possible
 - 2 bytes is common
 - Sneaky assembly code for operations on inverted files
 - Pack deltas into aligned 64 bit word
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

2/5/2013 2:21 PM

8

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



2/5/2013 2:21 PM

9

Processing Simple Queries

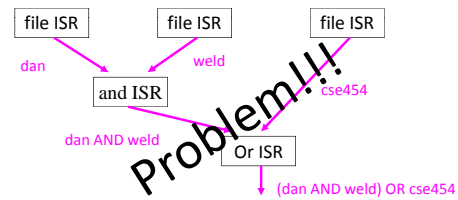
- User searches for “weld”
- Open ISR on “weld”
 - Uses hash table to avoid scanning entire file
- Next(), next(), next()
 - returns locations containing the word

2/5/2013 2:21 PM

10

Combining ISRs

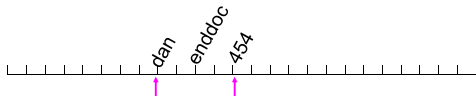
- And Compare locs on two streams
- Or Merges two or more ISRs
- Not Returns locations not in ISR (lazily)



2/5/2013 2:21 PM

11

What About File Boundaries?

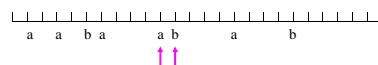


2/5/2013 2:21 PM

12

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”
 - $loc(A) \leq loc(B)$, $loc(B) \leq loc(A) + 1$



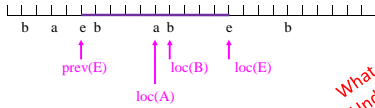
2/5/2013 2:21 PM

13

Two words on one page

- Let E be ISR for word `enddoc`
- Constraints for `a` AND `b`
 - $\text{prev}(E) \leq \text{loc}(A)$
 - $\text{loc}(A) \leq \text{loc}(E)$
 - $\text{prev}(E) \leq \text{loc}(B)$
 - $\text{loc}(B) \leq \text{loc}(E)$

$$\begin{aligned} \text{loc}(A) &\leq \text{loc}(B) + K \\ \text{prev}(A) &\leq \text{loc}(B) + K \\ \text{loc}(A) &\leq \text{prev}(B) + K \\ \text{prev}(A) &\leq \text{prev}(B) + K \end{aligned}$$



What if $\text{prev}(E)$ is Undefined?

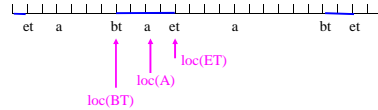
2/5/2013 2:21 PM

14

Advanced Search

- Query: `a` in Title of page
- Let BT, ET be ISRP of words
 - BeginTitle & EndTitle
- Constraints:
 - $\text{loc}(BT) \leq \text{loc}(A)$
 - $\text{loc}(A) \leq \text{loc}(ET)$

$$\begin{aligned} \text{loc}(A) &\leq \text{loc}(B) + K \\ \text{prev}(A) &\leq \text{loc}(B) + K \\ \text{loc}(A) &\leq \text{prev}(B) + K \\ \text{prev}(A) &\leq \text{prev}(B) + K \end{aligned}$$



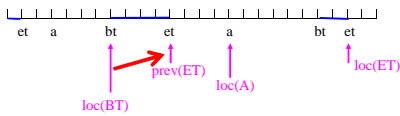
2/5/2013 2:21 PM

15

Advanced Search

- Query: `a` in Title of page
- Let BT, ET be ISRP of words
 - BeginTitle & EndTitle
- Constraints:
 - $\text{loc}(BT) \leq \text{loc}(A)$
 - $\text{loc}(A) \leq \text{loc}(ET)$
 - $\text{prev}(ET) \leq \text{loc}(BT)$ X

$$\begin{aligned} \text{loc}(A) &\leq \text{loc}(B) + K \\ \text{prev}(A) &\leq \text{loc}(B) + K \\ \text{loc}(A) &\leq \text{prev}(B) + K \\ \text{prev}(A) &\leq \text{prev}(B) + K \end{aligned}$$



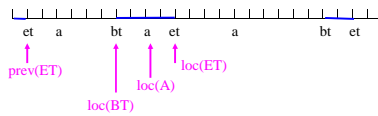
2/5/2013 2:21 PM

16

Advanced Search

- Query: `a` in Title of page
- Let BT, ET be ISRP of words
 - BeginTitle & EndTitle
- Constraints:
 - $\text{loc}(BT) \leq \text{loc}(A)$
 - $\text{loc}(A) \leq \text{loc}(ET)$
 - $\text{prev}(ET) \leq \text{loc}(BT)$

$$\begin{aligned} \text{loc}(A) &\leq \text{loc}(B) + K \\ \text{prev}(A) &\leq \text{loc}(B) + K \\ \text{loc}(A) &\leq \text{prev}(B) + K \\ \text{prev}(A) &\leq \text{prev}(B) + K \end{aligned}$$



2/5/2013 2:21 PM

17

Neat, Huh?

- Very flexible
- Eg "Who points to my page?"
- Are we done?

2/5/2013 2:21 PM

18

Implementing the Solver

Constraint Types

- $\text{loc}(A) \leq \text{loc}(B) + K$
- $\text{prev}(A) \leq \text{loc}(B) + K$
- $\text{loc}(A) \leq \text{prev}(B) + K$
- $\text{prev}(A) \leq \text{prev}(B) + K$

Remember: Index Stream Readers

ISR Methods

- $\text{loc}(\text{ISR})$ return current location
- $\text{next}(\text{ISR})$ advance to next location
- $\text{seek}(\text{ISR}, X)$ advance to next loc after X
- $\text{prev}(\text{ISR})$ return previous location

2/5/2013 2:21 PM

19

Solver Algorithm

loc(ISR)	return cur loc
next(ISR)	adv to nxt loc
seek(ISR, X)	return it
prev(ISR)	adv to nxt loc > X
	return it
	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)

2/5/2013 2:21 PM Copyright © 2000-2009 D.S.Weld 20

Solver Algorithm

loc(ISR)	return cur loc
next(ISR)	adv to nxt loc
seek(ISR, X)	return it
prev(ISR)	adv to nxt loc > X
	return it
	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)

2/5/2013 2:21 PM Copyright © 2000-2009 D.S.Weld 21

Solver Algorithm

loc(ISR)	return cur loc
next(ISR)	adv to nxt loc
seek(ISR, X)	return it
prev(ISR)	adv to nxt loc > X
	return it
	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:** $loc(A) \leq prev(B) + K$

2/5/2013 2:21 PM Copyright © 2000-2009 D.S.Weld 22

Solver Algorithm

loc(ISR)	return cur loc
next(ISR)	adv to nxt loc
seek(ISR, X)	return it
prev(ISR)	adv to nxt loc > X
	return it
	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),

2/5/2013 2:21 PM Copyright © 2000-2009 D.S.Weld 23

Solver Algorithm

loc(ISR)	return cur loc
next(ISR)	adv to nxt loc
seek(ISR, X)	return it
prev(ISR)	adv to nxt loc > X
	return it
	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)

2/5/2013 2:21 PM Copyright © 2000-2009 D.S.Weld 24

Solver Algorithm

loc(ISR)	return cur loc
next(ISR)	adv to nxt loc
seek(ISR, X)	return it
prev(ISR)	adv to nxt loc > X
	return it
	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)

Copyright © 2000-2009 D.S.Weld 25

Solver Algorithm

```
while (unsatisfied_constraints)
  satisfy_constraint(choose_unsat_constraint())
```

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

Heuristic??

Which choice advances a stream the furthest?

Copyright © 2000-2009 D.S.Weld

26

Update

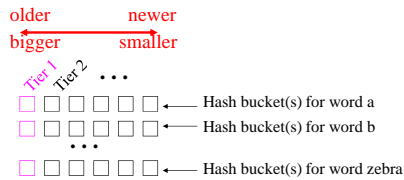
- **Can't insert in the middle of an inverted file**
- **Must rewrite the entire file**
 - Naïve approach: need space for two copies
 - Slow since file is huge
- **Split data along two dimensions**
 - **Buckets** solve disk space problem
 - **Tiers** alleviate small update problem

2/5/2013 2:21 PM

27

Buckets & Tiers

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

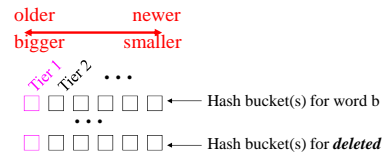


2/5/2013 2:21 PM

28

What if Word *Removed* from Doc?

- **Delete documents by adding deleted word**
-
- deleted
a
- **Expunge deletions when merging tier 1**



2/5/2013 2:21 PM

29