
**11-442 / 11-642:
Search Engines**

Introduction to Search

Jamie Callan
Carnegie Mellon University
callan@cs.cmu.edu

1

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean
- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

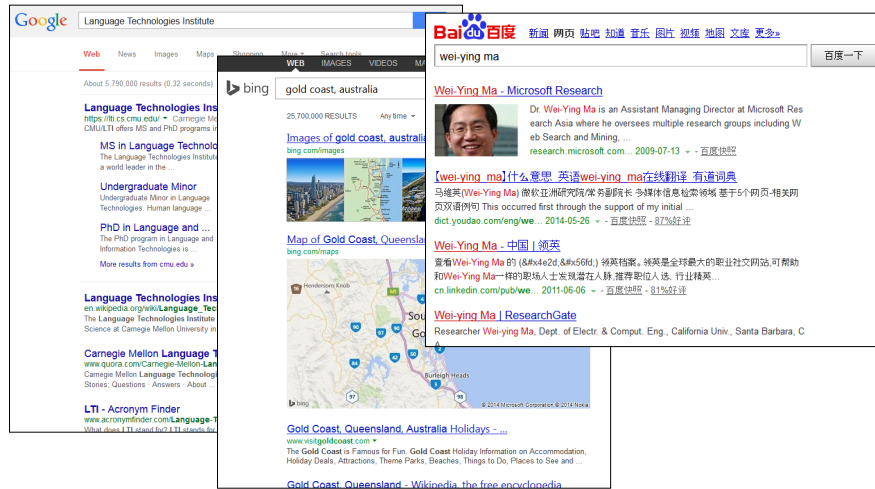
- Later lectures explore these topics in greater detail

2

© 2021, Jamie Callan

2

Probably You are an Experienced Search Engine User



3

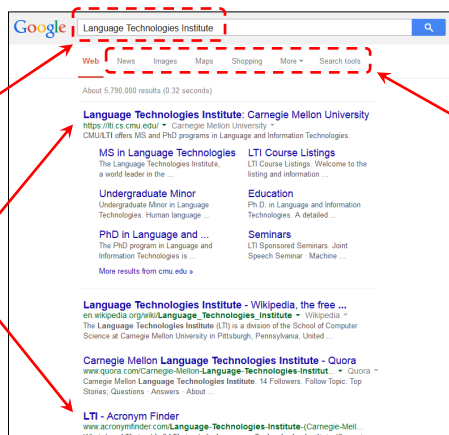
© 2021, Jamie Callan

3

How Does a Search Engine Developer or Researcher View the Search Process?

A query that expresses an information need

Documents (unstructured information)



Vertical search engines

4

© 2021, Jamie Callan

4

How Does a Search Engine Developer or Researcher View the Search Process?

A person starts with an information need

- The query is an approximate description of the information need

The person searches a corpus of unstructured information

- Documents

Goal: Find documents that satisfy the information need

- Search, retrieval

This lecture and the next present a simple end-to-end solution

- The “big picture”
- Later lectures go into more detail & more advanced material

5

© 2021, Jamie Callan

5

Simple End-to-End Solutions

Requirements

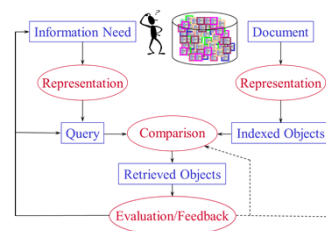
- A way of representing information needs
- A way of representing document content
- A comparison or matching process

Initial solutions

- Boolean queries
- Exact-match retrieval models (unranked and ranked)

These solutions are primitive, but they are still used today

- Fast, easy to build, easy to understand
- Often sufficient for a particular task



6

© 2021, Jamie Callan

6

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean
- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

- Later lectures explore these topics in greater detail

7

© 2021, Jamie Callan

7

Representing the Information Need



Exact-match retrieval models assume that a person can describe the information need as a Boolean query

- Relational database systems also make this assumption
- Most people are not good at creating Boolean queries
- Even well-trained people overestimate the quality of their queries



Examples:

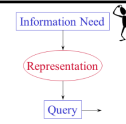
- Angelina AND Jolie
- (Angelina AND Jolie) OR (Brad AND Pitt)
- (Angelina NEAR/2 Jolie) OR (Brad NEAR/2 Pitt)
 - NEAR/n is similar to a phrase operator
 - Match if terms are in this order, separated by a distance $\leq n$

8

© 2021, Jamie Callan

8

Query Trees



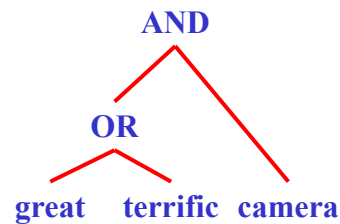
Query: (great OR terrific) AND camera

Search engines represent the query as a tree

- Nodes are query operators
- Leaves are index terms

Does query q match document d?

- Use depth-first evaluation
- More about this later....



9

© 2021, Jamie Callan

9

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean
- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

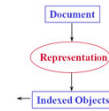
- Later lectures explore these topics in greater detail

10

© 2021, Jamie Callan

10

Representing the Document: An Example Document



A Great Choice.

Review by topjimmy5150

★★★★★ April, 21 2003

I have been looking and looking for a new camera to replace our bulky, but simple and reliable (but only fair picture taker) Sony Mavica FD73. My other choice (Besides the more expensive Nikon Coolpix 3100) was the (also more expensive) Sony Cybershot P72. I recommend any of these cameras, and I was set to buy the Sony, but at the last minute I cheaped out and bought the 2100. No regrets. I bought the camera (along with 128mb memory card (the stock 16mb card will be kept in the bag as a spare) and carrying case) at the new Best Buy in Harrisburg, PA. I also bought a set of 4 Nickle-Metal Hydride rechargable batteries and charger at Walmart for less than \$20. I keep 2 in the camera and two in the charger/in the camera bag along with the original Lithium battery pack as spares.

Hands down, the best feature of this camera is it's compact design. It is very small. My family likes to go camping during the summer, and last year we found the Mavica too

....

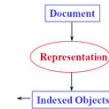
(topjimmy5150, Epinions.com)

11

© 2021, Jamie Callan

11

Representing the Document



How should the contents of a document be represented?

Today, assume that we will use words from the document

- Free-text indexing: Use just some of the words
 - Developed first, but ... which words?
- Full-text indexing: Use most or all of the words
 - Most search engines do this

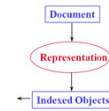
Later lectures consider other possibilities

12

© 2021, Jamie Callan

12

The Binary Full-Text Representation



Record which words occur in which documents (0: doesn't occur, 1: occurs)

- Invented first
- A tabular representation is simple (but very inefficient)

Vocabulary (Index Terms) |V|

	a	abba	abhor	ability	able	away	...	zoom
Doc₁	0	0	0	1	1	1	...	1
Doc₂	1	1	0	0	1	1	...	0
Doc_n	0	0	1	1	0	1	...	0

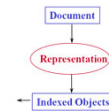
Corpus: Document collection

13

© 2021, Jamie Callan

13

The Binary Full-Text Representation (The Bag of Words)



In the binary full-text representation, position and frequency are ignored

- The document is a “bag of words”
 - Or other features (covered later)

This is a simple representation of meaning

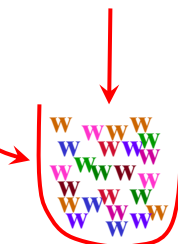
- But...surprisingly effective for search and other tasks (e.g., classification)

Full Review

I have been looking and looking for a new camera to replace our bulky, but simple and reliable (but only fair picture taken) Sony Mavica FD73. My other choice (Besides the more expensive Nikon Cooler 3100) was the (also more expensive) Sony Cybershot P72. I recommend any of these cameras, and I was set to buy the Sony, but at the last minute I cheaped out and bought the 2100. No regrets. I bought the camera (along with 128mb memory card the stock 16mb card will be used in the bag as a spare) and carrying case) at the new Best Buy in Harrisburg, PA. I also bought a set of 4 Nickel-Metal Hydride rechargeable batteries and charger at Walmart for less than \$20. I keep 2 in the camera and two in the charger in the camera bag along with the original Lithium battery pack as spares.

Hands down, the best feature of this camera is its compact design. It is very small. My family likes to go camping during the summer, and last year we found the Mavica too cumbersome to haul around. The 2100 is perfect size. It will easily slip into a shirt pocket. I like the way it looks and I like the way it feels in my hand. It feels perfect.

The photo quality is top notch in the 2mp range. I really wanted a 3mp camera, but



A bag of words

14

© 2021, Jamie Callan

14

Frequency-Based Full-Text Representation



Record the frequency of each word in each document

- More effective for search than the binary representation
- A tabular representation is simple (but very inefficient)

Vocabulary (Index Terms) |V|

	a	abba	abhor	ability	able	away	...	zoom
Doc₁	0	0	0	7	3	4	...	2
Doc₂	4	5	0	0	1	2	...	0
Doc_n	6	0	1	3	0	1	...	0

Corpus: Document collection

15

© 2021, Jamie Callan

15

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean
- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

- Later lectures explore these topics in greater detail

16

© 2021, Jamie Callan

16

Data Structures for Index Terms

Task: Evaluate the query “zoom AND away AND message”

One could compare to each document (row)

- Invented first
- Complexity is $O(|C| \times |V|)$

		Vocabulary (Index Terms) V							
		a	abba	abhor	ability	able	away	...	zoom
Corpus C	Doc ₁	0	0	0	1	1	1	...	1
	Doc ₂	1	1	0	0	1	1	...	0
	...	:	:	:	:	:	:	:	:
	Doc _n	0	0	1	1	0	1	...	0

Most terms are rare (occur in few documents)

- The vocabulary V is huge
- Nearly all documents fail to match the query
- Most of the $O(|C| \times |V|)$ effort is wasted effort

17

© 2021, Jamie Callan

17

Data Structures for Index Terms: Inverted Lists

		Vocabulary (Index Terms) V							
		a	abba	abhor	ability	able	away	...	zoom
Corpus C	Doc ₁	0	0	0	1	1	1	...	1
	Doc ₂	1	1	0	0	1	1	...	0
	...	:	:	:	:	:	:	:	:
	Doc _n	0	0	1	1	0	1	...	0

**Inverted
Lists**

Task: Evaluate the query “zoom AND away AND message”

One could compare to query terms (columns)

- Columns are called inverted lists
- Complexity is $O(|C| \times |Q|)$

		Vocabulary (Index Terms) V							
		a	abba	abhor	ability	able	away	...	zoom
Corpus C	Doc ₁	0	0	0	1	1	1	...	1
	Doc ₂	1	1	0	0	1	1	...	0
	...	:	:	:	:	:	:	:	:
	Doc _n	0	0	1	1	0	1	...	0

**Inverted
Lists
Really important
data structure!**

Most terms are rare (occur in few documents)

- Nearly all documents fail to match the query
- More efficient
 - ... but still, most of the $O(|C| \times |Q|)$ effort is wasted effort

18

© 2021, Jamie Callan

18

Are Fixed-Length Inverted Lists A Good Idea?

Vocabulary (Index Terms) - (t)									
Doc	1	2	3	4	5	6	7	8	9
Doc1	0	0	0	1	1	1	1	1	1
Doc2	1	1	0	0	1	1	1	1	0
Doc3	1	1	1	1	1	1	1	1	1
Doc4	0	0	1	1	0	1	1	1	0

Inverted Lists

Early search engines used fixed-length inverted lists

- Bit-vector operations are fast and easy to parallelize
- Very inefficient (1 bit or integer per document)

apple
1
0
0
1
:
0
1

Zipf's Law predicts that the median term occurs twice

- Rank of last term: $\text{Constant} / \text{Frequency} = (0.1 \times N) / 1$
- Rank of median term: $\frac{1}{2} \text{ rank of last term} = \frac{1}{2} \times (0.1 \times N) / 1$
- Frequency of median term: $\text{Constant} / \text{Rank of median term} = (0.1 \times N) / [\frac{1}{2} \times (0.1 \times N) / 1] = 1 / \frac{1}{2} = 2$

Thus, inverted lists use sparse representations

19

© 2021, Jamie Callan

19

Sparse Representation of Inverted Lists

Simple approach: Store ids of documents that contain the word

- E.g., apple: length=18, docids: 1, 5, 6, 9, ...
 - The term 'apple' occurs in 18 documents

A more typical notation (used often in this course)

- $df_t=18$, docids 1, 5, 6, 9, ...
 - df_t : document frequency (number of documents containing term t)
 - docids: document identifiers

You must know this data structure & more advanced variants

20

© 2021, Jamie Callan

20

Different Types of Inverted Lists for the Term “apple”

Vocabulary (Index Term) (V)					
	a	abba	abba	abba	apple
Doc1	0	0	0	1	1
Doc2	1	1	0	0	1
Doc3	1	1	1	1	1
Doc4	0	0	1	1	0

Corpus (C)

Inverted Lists

Different types of inverted lists support different capabilities

Binary Inverted lists	Supports
df: 4356	• Unranked
docid: 42	• Boolean
docid: 94	• AND, OR
:	• SYN

Frequency Inverted lists	Also supports
df: 4356	• Ranking
docid: 42	• SUM
tf: 3	• WSUM
docid: 94	
tf: 1	
:	

Positional Inverted lists	Also supports
df: 4356	• NEAR/n
docid: 42	
tf: 3	
locations: 14	
83	
157	
docid: 94	
tf: 1	
locations: 65	
:	

df: document frequency
docid: sequential document ids
tf: term frequency ($tf_{t,d}$)
locs: locations where t appears in d

21

© 2021, Jamie Callan

21

Inverted Indexes

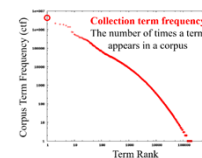
Vocabulary (Index Term) (V)					
	a	abba	abba	abba	apple
Doc1	0	0	0	1	1
Doc2	1	1	0	0	1
Doc3	1	1	1	1	1
Doc4	0	0	1	1	0

Corpus (C)

Inverted Lists

After indexing, there are many inverted lists

- One per term in the vocabulary (typically 10^6 to 10^8)
- Very skewed size distribution (Zipf's Law)
- Very skewed access patterns



An inverted index consists of two parts

- Inverted file(s) that contain inverted lists
 - An object database containing the inverted lists
- An access mechanism
 - Term string → inverted list
 - Term id → inverted list
 - Sometimes combined with the term dictionary

22

© 2021, Jamie Callan

22

Inverted Indexes: Two Common Access Methods

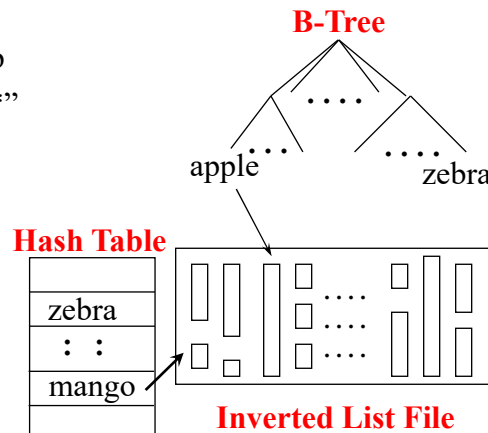
How is a file of inverted lists accessed?

- **B-Tree (B+ Tree, B* Tree, etc)**

- Exact-match and range-based lookup
 - » “apple”, “apple – apples”, “appl*”
- $O(\log n)$ lookups to find a list
- Usually easy to expand

- **Hash table**

- Exact-match lookup
 - » “apple”
- $O(1)$ lookups to find a list
- May be complex to expand



23

© 2021, Jamie Callan

23

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean
- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

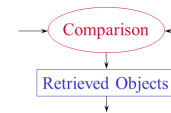
- Later lectures explore these topics in greater detail

24

© 2021, Jamie Callan

24

Retrieval Model #1: Unranked Boolean



Model: Retrieve documents iff they satisfy a Boolean expression

- **Examples:** “michelle AND obama”, “biden OR trump”
- The query specifies exact relevance criteria
 - “Exact match” retrieval
- The set of matching (“retrieved”) documents is unordered
 - Often sorted by date

Query operators:

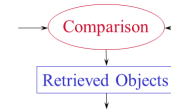
- AND, OR, ANDNOT, NEAR, WINDOW, DATE, BEFORE, ...
- Typically these systems have rich query languages

25

© 2021, Jamie Callan

25

Retrieval Model #1: Unranked Boolean



**This approach to document retrieval was invented first
...and was the dominant model until the early 1990s
...but it is no longer state-of-the-art**

Why?

- Most people find it difficult to construct good Boolean queries
- Documents are returned in no particular order

However, it is still used in many systems, and still important

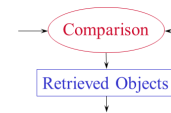
- E.g., WestLaw, PubMed, first pass in Web search engines, ...

26

© 2021, Jamie Callan

26

Retrieval Model #2: Ranked Boolean



Model: Retrieve documents iff they satisfy a Boolean expression

- Order the matching (“retrieved”) documents by scores
- Document scores can be anything you want
 - But, typically Ranked Boolean systems have simple scores

Note: Unranked Boolean systems have implied document scores

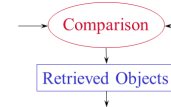
- 1: Query matches document
- 0: Query doesn’t match document

27

© 2021, Jamie Callan

27

Ranked Boolean: Calculating Scores



What is the score when a query term t occurs in document d ?

$tf_{t,d}$

- The frequency of term t in document d
- Invented first
- Easy to implement
- Only considers document d

$tf_{t,d} \times idf_t = tf_{t,d} \times \log(N / df_t)$

- N : # of documents in the corpus
- df_t : # of documents that contain t
- Reward frequent terms in d
- Penalize frequent terms in the corpus
- More effective

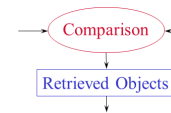
Both types of weights are used

28

© 2021, Jamie Callan

28

Ranked Boolean: Calculating Scores



Boolean queries have operators such as AND and OR

- cat AND mouse; john AND paul AND george AND ringo
- rich OR poor; obama OR bush OR clinton OR reagan

A prefix representation makes the query structure more apparent

- AND (cat mouse) AND (john paul george ringo)
- OR (rich poor) OR (obama bush clinton reagan)

Notation: $q_{\text{operator}}(q_1 \dots q_n)$

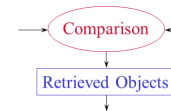
- $q_{\text{AND}}(\text{cat mouse})$ $q_{\text{AND}}(\text{john paul george ringo})$
- $q_{\text{OR}}(\text{rich poor})$ $q_{\text{OR}}(\text{obama bush clinton reagan})$

29

© 2021, Jamie Callan

29

Ranked Boolean: Calculating Scores



What is the score for AND operator $q_{\text{AND}}(q_1 \dots q_n)$ on document j ?

- $\text{score}(q_{\text{AND}}(q_1 \dots q_n), d_j) = \text{MIN}(\text{score}(q_1, d_j), \dots, \text{score}(q_n, d_j))$

What is the score for OR operator $q_{\text{OR}}(q_1 \dots q_n)$ on document j ?

- $\text{score}(q_{\text{OR}}(q_1 \dots q_n), d_j) = \text{MAX}(\text{score}(q_1, d_j), \dots, \text{score}(q_n, d_j))$
 - Consistent with the AND operator
- $\text{score}(q_{\text{OR}}(q_1 \dots q_n), d_j) = \text{MEAN}(\text{score}(q_1, d_j), \dots, \text{score}(q_n, d_j))$
 - Rewards documents that match many query terms
 - The semantics of an OR operator do not require this behavior
 - But, it is the behavior that people expect
 - Typically a little more effective than MAX

30

© 2021, Jamie Callan

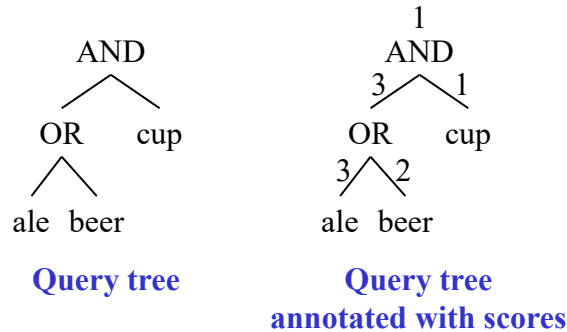
30

Ranked Boolean: Calculating Scores

Document:

... ale ... cup ... ale ... beer ... ale ... beer

Query: (ale OR beer) AND cup

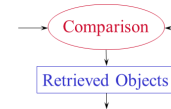


31

© 2021, Jamie Callan

31

Retrieval Model #2: Ranked Boolean



Advantages

- Very efficient
- Predictable, easy to explain, structured queries
- Works well enough when searchers know exactly what is wanted
- Results ordered by how redundantly a document satisfies a query
- Other term weighting methods can be used, too

Disadvantages

- It's still an Exact-Match model
- Usually it is difficult to get a good balance of Precision and Recall

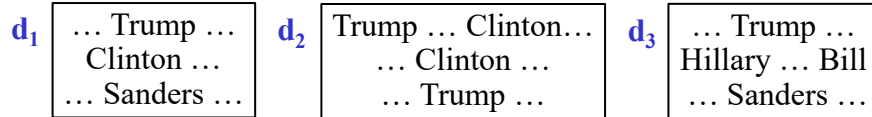
32

© 2021, Jamie Callan

32

Exact-Match Retrieval: Unranked vs. Ranked Boolean Retrieval

Query: Trump AND Clinton



Three retrieval methods

Unranked Boolean

d_1
 d_2
(arbitrary order)

Ranked Boolean

d_2
 d_1

Best Match

d_2
 d_1
 d_3

Which ranking is best?

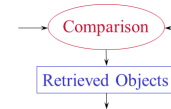
- It depends on the task ... sometimes unranked Boolean is enough

33

© 2021, Jamie Callan

33

Are Exact-Match Models Still Relevant?



Many people prefer exact-match Boolean models

- Professional searchers (e.g., librarians, paralegals)
- Some Web surfers (e.g., “Advanced Search” feature)
- What do they like? Control, predictability, understandability
- Preferred by 70% of WESTLAW searchers in a 2007 survey

-- James Allan, 2007

Exact-match Boolean is a low-level part of Web search engines

- Massive corpus makes efficiency important
- Massive corpus makes partial matching less important

34

© 2021, Jamie Callan

34

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean
- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

- Later lectures explore these topics in greater detail

35

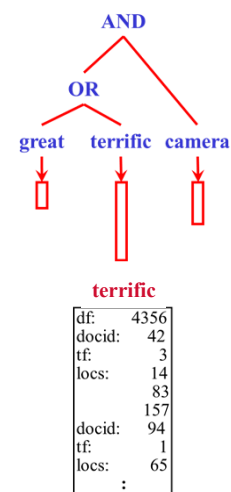
© 2021, Jamie Callan

35

Document Retrieval

There are three approaches to query processing

- Term-at-a-Time (TAAT)
- Document-at-a-Time (DAAT)
- TAAT / DAAT hybrids
 - Important in large-scale systems, but not covered in this class



36

© 2021, Jamie Callan

36

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Key ideas

- Fully process $list_i$ before proceeding to $list_{i+1}$
- Each time a list is processed, partial document scores are updated

37

© 2021, Jamie Callan

37

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: (a AND b) OR c

38

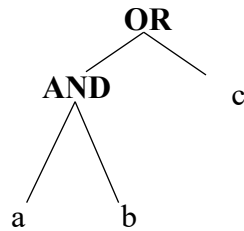
© 2021, Jamie Callan

38

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: (a AND b) OR c

Parsed
Query



39

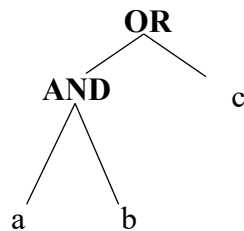
© 2021, Jamie Callan

39

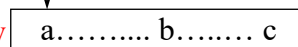
Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: (a AND b) OR c

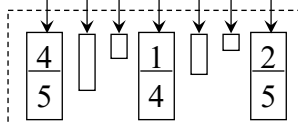
Parsed
Query



Term
Dictionary



Inverted
List File



40

© 2021, Jamie Callan

40

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

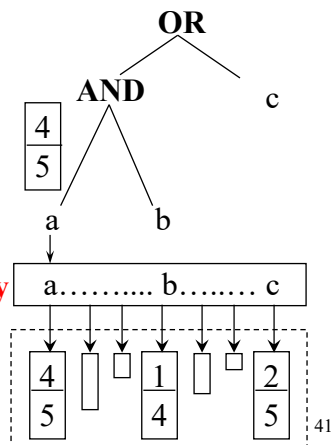
**Parsed
Query**

Query: (a AND b) OR c

1. Read inverted list for 'a'
from inverted list database

**Term
Dictionary**

**Inverted
List File**



41

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

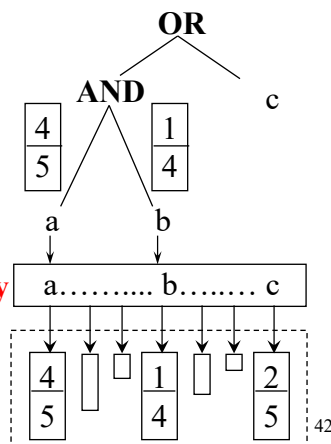
**Parsed
Query**

Query: (a AND b) OR c

1. Read inverted list for 'a'
from inverted list database
2. Read inverted list for 'b'
from inverted list database

**Term
Dictionary**

**Inverted
List File**



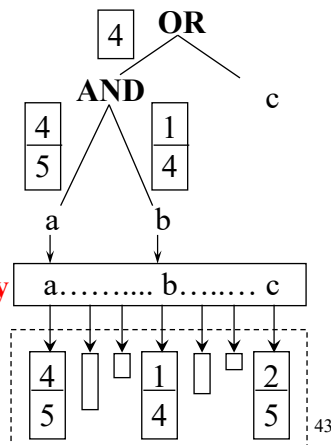
42

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

**Parsed
Query**

**Term
Dictionary**

**Inverted
List File**



Query: $(a \text{ AND } b) \text{ OR } c$

1. Read inverted list for 'a' from inverted list database
2. Read inverted list for 'b' from inverted list database
3. AND operator: Intersect the inverted lists for 'a' and 'b'

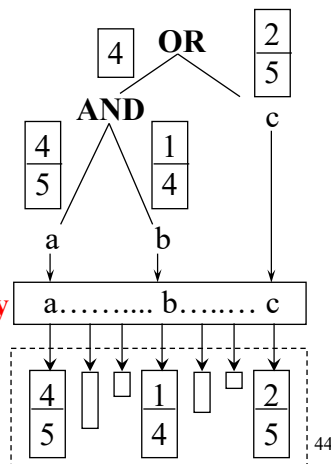
43

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

**Parsed
Query**

**Term
Dictionary**

**Inverted
List File**

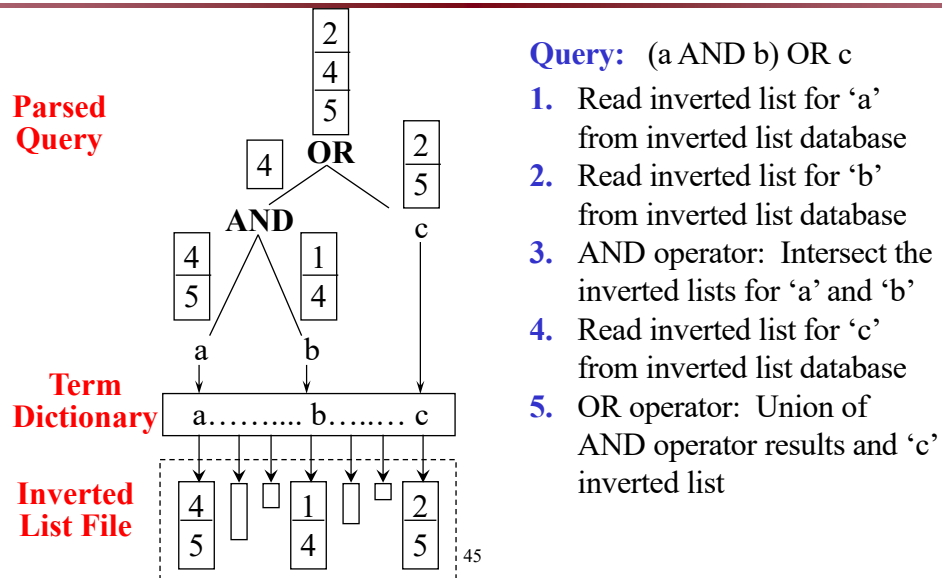


Query: $(a \text{ AND } b) \text{ OR } c$

1. Read inverted list for 'a' from inverted list database
2. Read inverted list for 'b' from inverted list database
3. AND operator: Intersect the inverted lists for 'a' and 'b'
4. Read inverted list for 'c' from inverted list database

44

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation



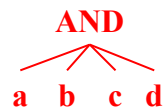
45

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b c d)

Evaluation strategy

- Retrieve a
- Retrieve b
- a AND b \rightarrow Result_{AND_1}
- Retrieve c
- Result_{AND_1} AND c \rightarrow Result_{AND_2}
- Retrieve d
- Result_{AND_2} AND d \rightarrow Result_Q



46

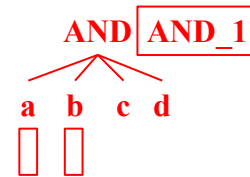
46

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b c d)

Evaluation strategy

- Retrieve a
- Retrieve b
- $a \text{ AND } b \rightarrow \text{Result}_{\text{AND_1}}$
- Retrieve c
- $\text{Result}_{\text{AND_1}} \text{ AND } c \rightarrow \text{Result}_{\text{AND_2}}$
- Retrieve d
- $\text{Result}_{\text{AND_2}} \text{ AND } d \rightarrow \text{Result}_Q$



47

© 2021, Jamie Callan

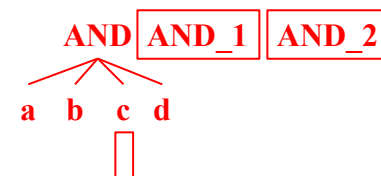
47

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b c d)

Evaluation strategy

- Retrieve a
- Retrieve b
- $a \text{ AND } b \rightarrow \text{Result}_{\text{AND_1}}$
- Retrieve c
- $\text{Result}_{\text{AND_1}} \text{ AND } c \rightarrow \text{Result}_{\text{AND_2}}$
- Retrieve d
- $\text{Result}_{\text{AND_2}} \text{ AND } d \rightarrow \text{Result}_Q$



48

© 2021, Jamie Callan

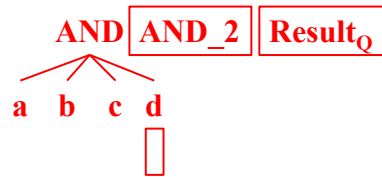
48

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b c d)

Evaluation strategy

- Retrieve a
- Retrieve b
- $a \text{ AND } b \rightarrow \text{Result}_{\text{AND_1}}$
- Retrieve c
- $\text{Result}_{\text{AND_1}} \text{ AND } c \rightarrow \text{Result}_{\text{AND_2}}$
- Retrieve d
- $\text{Result}_{\text{AND_2}} \text{ AND } d \rightarrow \text{Result}_Q$



49

© 2021, Jamie Callan

49

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b c d)

Evaluation strategy

- Retrieve a
- Retrieve b
- $a \text{ AND } b \rightarrow \text{Result}_{\text{AND_1}}$
- Retrieve c
- $\text{Result}_{\text{AND_1}} \text{ AND } c \rightarrow \text{Result}_{\text{AND_2}}$
- Retrieve d
- $\text{Result}_{\text{AND_2}} \text{ AND } d \rightarrow \text{Result}_Q$



50

© 2021, Jamie Callan

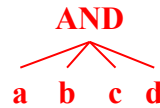
50

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b c d)

Characteristics

- Each query operator stores in RAM up to 3 lists simultaneously
 - arg_1 , arg_2 , result
- Peak query operator memory usage for this query
 - 3 lists in RAM simultaneously
 - $\text{size}(\text{arg}_1) + \text{size}(\text{arg}_2) + \text{size}(\text{result})$ bytes



51

© 2021, Jamie Callan

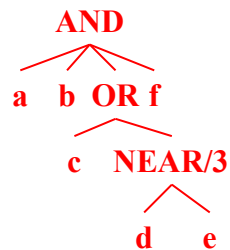
51

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Query: #AND (a b #OR (c #NEAR/3 (d e)) f)

Peak memory usage (probably)

- 5 lists in memory simultaneously
- $\text{size}(a \text{ AND } b) + \text{size}(c) + \text{size}(d) + \text{size}(e) + \text{size}(d \text{ NEAR/3 } e)$ bytes



52

© 2021, Jamie Callan

52

Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

Easy to understand and build

- Thus, we cover them first

Very efficient

- Little wasted effort
- This is more apparent when we consider DAAT

Memory usage is uncontrolled

- A query of depth d must store $d+2$ lists in RAM

Causes of memory problems

- Queries with frequent terms (long lists)
- Complex queries (more lists)
- Systems that process multiple queries in parallel (contention)

Rarely used in large systems

53

© 2021, Jamie Callan

53

Two Lecture Outline

A quick introduction to...

- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Indexes
 - Inverted lists
- Exact match retrieval
 - Unranked Boolean
 - Ranked Boolean

- Document retrieval
 - TAAT
 - DAAT
- Query operators
 - Types of query operators
 - The NEAR operator

Goal: Provide an overview of search (“the Big Picture”)

- Later lectures explore these topics in greater detail

54

© 2021, Jamie Callan

54