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

# Introduction to Search

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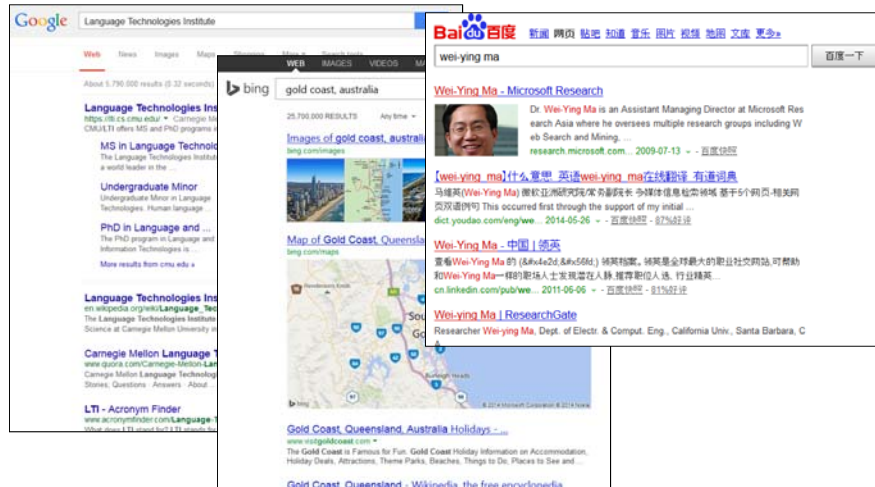
---

# Two Lecture Outline

**A quick introduction to…**
- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Exact match retrieval
    - Unranked Boolean
    - Ranked Boolean
- Indexing
    - Inverted lists
    - Term dictionary
- Query processing
    - TAAT
    - DAAT
- Query operators

**Goal:** Provide an overview of search ("the Big Picture")
- Later lectures explore these topics in greater detail

2

## Probably You are an Experienced Search Engine User
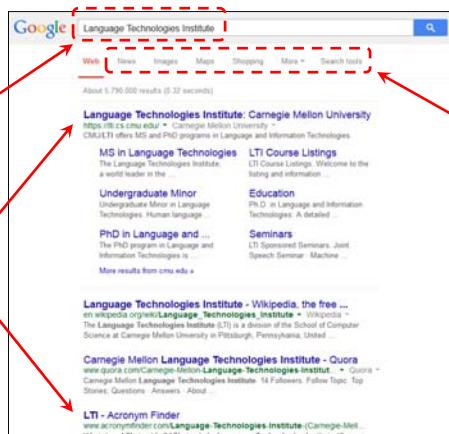
3                                                          © 2017, Jamie Callan

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

**A query that expresses an information need**

**Vertical search engines**

**Documents (unstructured information)**

4                                                          © 2017, Jamie Callan

## 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

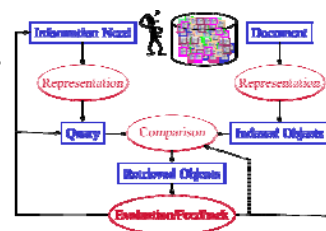## Simple End-to-End Solutions

**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

**Two Lecture Outline**

**A quick introduction to…**
- Ad-hoc retrieval
- **Information needs & queries**
- Document representation
- Exact match retrieval
  - Unranked Boolean
  - Ranked Boolean

- Indexing
  - Inverted lists
  - Term dictionary
- Query processing
  - TAAT
  - DAAT
- Query operators

**Goal:** Provide an overview of search ("the Big Picture")
- Later lectures explore these topics in greater detail

7 © 2017, Jamie Callan

---

**Representing the Information Need**

**Exact-match retrieval models assume that a person can describe the information need as a <u>Boolean query</u>**
- 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
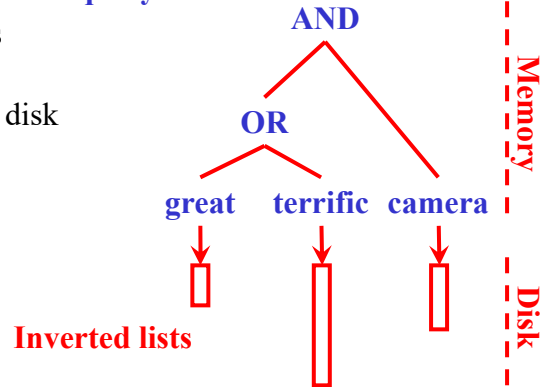
8 © 2017, Jamie Callan

## Query Trees

**Query:** (great OR terrific) AND camera

**Search engines represent the query as a tree**
- Leaves are index terms
- Leaves are pointers to inverted lists stored on disk

AND

OR

great    terrific  camera

Memory

Disk

**Inverted lists**

9

---

## Two Lecture Outline

**A quick introduction to…**
- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Exact match retrieval
  - Unranked Boolean
  - Ranked Boolean

- Indexing
  - Inverted lists
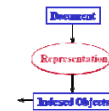  - Term dictionary
- Query processing
  - TAAT
  - DAAT
- Query operators

**Goal:** Provide an overview of search ("the Big Picture")
- Later lectures explore these topics in greater detail

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                                                                 © 2017, Jamie Callan

---

## Representing the Document

**How should the contents of a document be represented?**

**Today, assume that we will use words from the document**
- <u>Free-text indexing</u>:  Use just some of the words
    – Developed first, but … which words?
- <u>Full-text indexing</u>:  Use most or all of the words
    – Most search engines do this

**Later lectures consider other possibilities**

12                                                                 © 2017, Jamie Callan

## The Binary Full-Text Representation

**Record which words occur in which documents**
- Invented first
- Almost never used anymore, but a useful simplification

**Corpus Vocabulary   |V|**

|  | a | abba | abend | ability | able | about | … | zooms |
|---|---|---|---|---|---|---|---|---|
| **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 |

**Corpus |C|**

13

© 2017, Jamie Callan

---

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

**In the binary full-text representation, <u>position</u> and <u>frequency</u> are ignored**
- The document is a "<u>bag of words</u>"
  - Or other features (covered later)
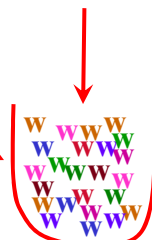
**This is a simple representation of meaning**
- But…surprisingly effective for search and other tasks (e.g., classification)

Full Review

I have been looking and looking for a new camera to replace our bulky, but simple and reliable (but only fair picture 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 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

14

© 2017, Jamie Callan

Page 7

## Frequency-Based Full-Text Representation

**Record the frequency of each word in each document**
- More effective for search than the binary representation
- The tabular representation is a useful simplification

**Corpus Vocabulary    |V|**

| | a | abba | abend | ability | able | about | … | zooms |
|---|---|---|---|---|---|---|---|---|
| **Doc₁** | 0 | 0 | 0 | 7 | 3 | 4 | … | 2 |
| **Doc₂** | 4 | 5 | 0 | 0 | 1 | 2 | … | 0 |
| ⁞ ⁞ ⁞ ⁞ | ⁞ | ⁞ ⁞ | ⁞ ⁞ | ⁞ ⁞ | ⁞ ⁞ | ⁞ ⁞ ⁞ | | ⁞ ⁞ ⁞ |
| **Docₙ** | 6 | 0 | 1 | 3 | 0 | 1 | … | 0 |

**Corpus |C|**

15

© 2017, Jamie Callan

---

## Two Lecture Outline

**A quick introduction to…**
- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Exact match retrieval
    - Unranked Boolean
    - Ranked Boolean
- Indexing
    - Inverted lists
    - Term dictionary
- Query processing
    - TAAT
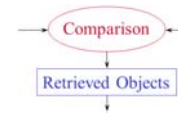    - DAAT
- Query operators

**Goal:** Provide an overview of search ("the Big Picture")
- Later lectures explore these topics in greater detail

16

© 2017, Jamie Callan

Page 8

## Retrieval Model #1:
## Unranked Boolean

**Model:** Retrieve documents <u>iff</u> they satisfy a Boolean expression
- **Examples:** "michelle AND obama", "clinton OR trump"
- The query specifies exact relevance criteria
    - <u>Exact match</u> retrieval
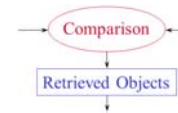- The matching set of retrieved documents is <u>unordered</u>
    - Often sorted by date

**Query operators:**
- AND, OR, ANDNOT, NEAR, DATE, …
- Typically these systems have sophisticated query languages
    - A rich language to describe the set of relevant documents

---

## 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?**
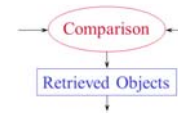- Documents are returned in no particular order

**However, it is still used in many systems, and still important**
- E.g., WestLaw, PubMed, first pass in Web search engines, …

**Retrieval Model #2:**
**Ranked Boolean**
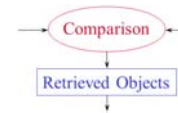
Comparison
Retrieved Objects

**Ranked Boolean addresses a flaw of Unranked Boolean**

- Model & operators are the same as for Unranked Boolean
- But … matching documents are ranked by a <u>document score</u>
  - Unranked Boolean document scores are 0 and 1
  - Ranked Boolean document scores can be anything
    - » Typical Ranked Boolean systems have simple scores

19 © 2017, Jamie Callan

---

**Ranked Boolean:**
**Calculating Scores**

Comparison
Retrieved Objects

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

- $\text{tf}_{t,d}$: The frequency of term $t$ in document $d$
- $\text{tf}_{t,d} \times \text{idf}_t = \text{tf}_{t,d} \times \log(N / \text{df}_t)$: Penalize common terms
  - N: Number of documents in the corpus
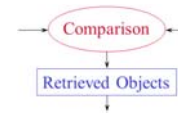  - $\text{df}_t$: The number of documents that contain term $t$

**Both types of weights are used**

- $\text{tf}_{t,d}$: Invented first, easy to implement, only considers the document
- $\text{tf}_{t,d} \times \text{idf}_t$: More effective
  - Rewards terms that are frequent in this document, but not frequent in the corpus

20 © 2017, Jamie Callan

Page 10

**Ranked Boolean:**
**Calculating Scores**

Comparison ←
Retrieved Objects

**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)
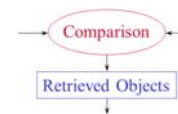- OR (rich poor)            OR (obama bush clinton reagan)

**Notation:**  $q_{operator} (q_1 \dots q_n)$
- $q_{AND}$ (cat mouse)          $q_{AND}$ (john paul george ringo)
- $q_{OR}$ (rich poor)          $q_{OR}$ (obama bush clinton reagan)

---

**Ranked Boolean:**
**Calculating Scores**

Comparison ←
Retrieved Objects

**What is the score for AND operator $q_{AND} (q_1 \dots q_n)$ on document j?**
- score $(q_{AND} (q_1 \dots q_n), d_j)$ = MIN (score $(q_1, d_j)$, …, score $(q_n, d_j)$)

**What is the score for OR operator $q_{OR} (q_1 \dots q_n)$ on document j?**
- score $(q_{OR} (q_1 \dots q_n), d_j)$ = MAX (score $(q_1, d_j)$, …, score $(q_n, d_j)$)
  – Consistent with the AND operator
- score $(q_{OR} (q_1 \dots q_n), d_j)$ = MEAN (score $(q_1, d_j)$, …, 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

## Retrieval Model #2:
## Ranked Boolean

**Comparison**

Retrieved Objects

**Advantages**

- <u>Very</u> 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

---

## Exact-Match Retrieval:
## Unranked vs. Ranked Boolean Retrieval

**Query:** Trump AND Clinton

$D_1$
| … Trump … |
| Clinton … |
| … Sanders … |

$D_2$
| Trump … Clinton… |
| … Clinton … |
| … Trump … |

$D_3$
| … Trump … |
| Hillary … Bill |
| … Sanders … |

**Three retrieval methods**

| Unranked Boolean | Ranked Boolean | Best Match |
|---|---|---|
| $D_1$ | $D_2$ | $D_2$ |
| $D_2$ | $D_1$ | $D_1$ |
| (arbitrary order) | | $D_3$ |

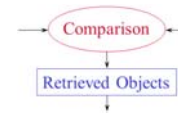**Which ranking is best?**

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

## 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

---

## Two Lecture Outline

**A quick introduction to…**
- Ad-hoc retrieval
- Information needs & queries
- Document representation
- Exact match retrieval
  - Unranked Boolean
  - Ranked Boolean

- Indexing
  - Inverted lists
    - Term dictionary
- Query processing
  - TAAT
  - DAAT
- Query operators

**Goal:**  Provide an overview of search ("the Big Picture")
- Later lectures explore these topics in greater detail

---

## Data Structures for Index Terms

**Search engines use <u>indexes</u> to process queries quickly**
- Indexes are data structures that allow things to be found quickly
- Present a key, get back an object

**What kinds of lookup does a search engine need?**
- Term $\rightarrow$ DocId +       "apple" $\rightarrow$ $doc_1$ $doc_{14}$ $doc_{32}$ …
- Attribute $\rightarrow$ Docid +     2014-08-06 $\rightarrow$ $doc_8$ $doc_{19}$ $doc_{63}$ …
- Term $\rightarrow$ Corpus statistics    "apple" $\rightarrow$ df=3,731; ctf=8,839; …
- Doc id $\rightarrow$ Attribute        $doc_{19}$ $\rightarrow$ 2014-01-04
- Doc id $\rightarrow$ Term +         $doc_{19}$ $\rightarrow$ "listen" "live" "WYEP" …
- Doc id $\rightarrow$ Document      $doc_{19}$ $\rightarrow$ 
- …

27        © 2017, Jamie Callan

---

## Data Structures for Index Terms

**Task:** Evaluate the query "ability AND about"

**One could compare to each document (row)**
- Invented first
- Rows are bit vectors
- Complexity is $O(|C| \times |V|)$
- Is this good?

|  | Corpus Vocabulary |V| | | | | | | |
|---|---|---|---|---|---|---|---|---|
|  | a | abba | abend | ability | able | about | … | zooms |
| Doc₁ | 1 | 0 | 0 | 1 | 1 | 1 | … | 1 |
| Doc₂ | 1 | 1 | 0 | 0 | 1 | 1 | … | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | | ⋮ |
| Docₙ | 1 | 0 | 1 | 1 | 0 | 1 | … | 0 |

Corpus |C|

**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

28        © 2017, Jamie Callan

Page 14

## Data Structures for Index Terms: Inverted Lists

**Task:** Evaluate the query "ability AND about"

**One could compare to query terms (columns)**

- Columns are bit vectors
- Columns are called <u>inverted lists</u>
- Complexity is $O(|C| \times |Q|)$

| | a | abba | abend | bility | able | bout | ... | zooms |
|---|---|---|---|---|---|---|---|---|
| Doc₁ | 1 | 0 | 0 | 1 | 1 | 1 | ... | 1 |
| Doc₂ | 1 | 1 | 0 | 0 | 1 | 1 | ... | 0 |
| :::: | : | : : | : : | : : | : : | : : : | | : : : |
| Docₙ | 1 | 0 | 1 | 1 | 0 | 1 | ... | 0 |

Corpus Vocabulary

Corpus

**Inverted Lists**

**<span style="color:red">Really important data structure!</span>**

**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

29 © 2017, Jamie Callan

---

## Are Fixed-Length Inverted Lists A Good Idea?

**Fixed-length inverted lists were used in early search engines**

- Simple to manage
- Bit-vector operations are fast and easy to parallelize
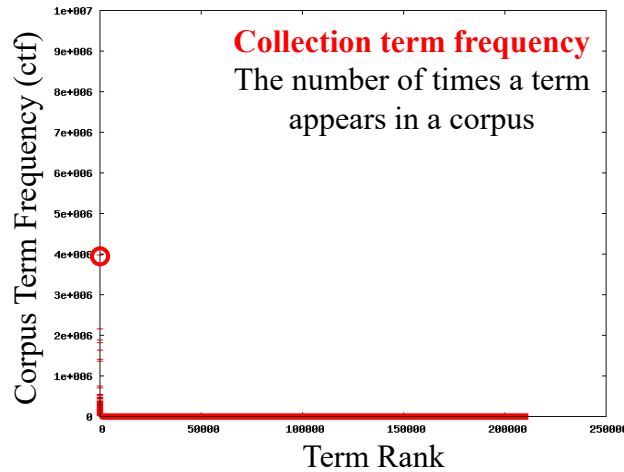
**Fixed-length inverted lists are inefficient but … how inefficient?**

- Length: $|C|$ bits
  - The number of documents in the corpus
- Number of bits set to 1 in a typical inverted list: $df_{typical}$
  - The <u>document frequency</u> of the term with the median term frequency
- How does $|C|$ compare to $df_{typical}$?

30 © 2017, Jamie Callan

## Term Frequency in the WSJ

**Collection term frequency**
The number of times a term appears in a corpus

- Y-axis: Corpus Term Frequency (ctf): 1e+007, 9e+006, 8e+006, 7e+006, 6e+006, 5e+006, 4e+006, 3e+006, 2e+006, 1e+006, 0
- X-axis: Term Rank: 0, 50000, 100000, 150000, 200000, 250000

**Wall Street Journal (1987-1992)**

| | |
|---|---|
| Documents: | 174K |
| Tokens: | 69M |
| Terms (types): | 211K |
| Megabytes: | 533 |

**Tokens:**
Word occurrences

**Types:**
Unique words

31

---

## Term Frequency in the WSJ

**Collection term frequency**
The number of times a term appears in a corpus

- Y-axis: Corpus Term Frequency (ctf): 1e+007, 1e+006, 100000, 10000, 1000, 100, 10, 1
- X-axis: Term Rank: 1, 10, 100, 1000, 10000, 100000, 1e+006

**Half of the vocabulary**

**Wall Street Journal (1987-1992)**

| | |
|---|---|
| Documents: | 174K |
| Tokens: | 69M |
| Terms (types): | 211K |
| Megabytes: | 533 |

**Tokens:**
Word occurrences

**Types:**
Unique words

32

Page 16

*16*

## Are Fixed-Length Inverted Lists A Good Idea?

**In ordinary text, $|C| \gg |df_{median}|$**
- In the Wall Street Journal, $|C| = 174K$ and $df_{median} = 2$
    - i.e., the median term appears in just <u>two</u> documents

**So usually fixed-length inverted lists are <u>very</u> inefficient**
- Often used for illustration
- Rarely used in real systems

33 © 2017, Jamie Callan

---

## Sparse Representation of Inverted Lists

**Simple approach:** Store ids of documents that contain the word
- **Full inverted list:** 10001100100…
    - The $i^{th}$ bit indicates whether the term occurs in the $i^{th}$ document
- **Sparse inverted list:** length=18, docids: 1, 5, 6, 9, …
    - The term occurs in 18 documents

**Usually in this course inverted lists are represented as**
- **$df_t$=18, docids 1, 5, 6, 9, …**
    **$df_t$:** document frequency (number of documents containing t)
    **docids:** document identifiers

**You <u>must</u> know this data structure & more advanced variants**

34 © 2017, Jamie Callan

**Inverted Lists**

**Different types of inverted lists support different capabilities**

- **Binary inverted lists**
  - Unranked retrieval
  - AND, OR, AND-NOT, FIELD Boolean operators
- **Frequency inverted lists**
  - Ranked retrieval
  - The above operators plus SUM, SYNONYM
- **Positional inverted lists**
  - Ranked retrieval
  - The above operators plus positional operators
    - » NEAR/n, SENTENCE/n, PASSAGE/n, WINDOW/n

---

**Different Types of Inverted Lists
for the Term "Apple"**

| Binary | Frequency | Positional | Positional, With Embedded Document Info |
|---|---|---|---|
| df: 4356<br>docid: 42<br>docid: 94<br>: | df: 4356<br>docid: 42<br>tf: 3<br>docid: 94<br>tf: 1<br>: | df: 4356<br>docid: 42<br>tf: 3<br>locs: 14<br>83<br>157<br>docid: 94<br>tf: 1<br>locs: 65<br>: | df: 4356<br>docid: 42<br>doclen: 357<br>tf: 3<br>locs: 14<br>83<br>157<br>docid: 94<br>doclen: 172<br>tf: 1<br>locs: 65<br>: |

df: document frequency
tf: term frequency ($tf_{t,d}$)
doclen: document length (in terms)

## Inverted Indexes

**After indexing, there are many <u>inverted lists</u>**
- 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 <u>index</u> consists of two parts**
- Inverted <u>file(s)</u> that contain inverted <u>lists</u>
    - An object database containing the inverted lists
- An <u>access mechanism</u>
    - Term string → inverted list
    - Term id → inverted list
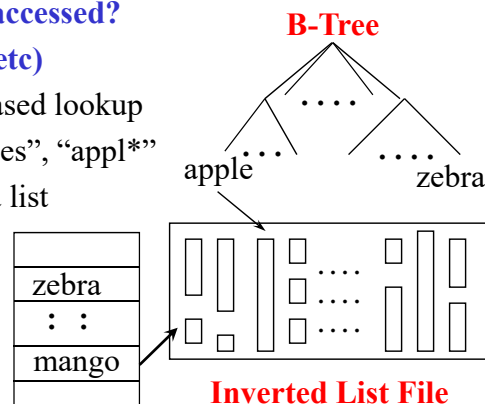    - Sometimes combined with the term dictionary

37

© 2017, Jamie Callan

---

## 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

**B-Tree**

apple          zebra

zebra
: :
mango

**Hash Table**

**Inverted List File**

38

© 2017, Jamie Callan

Page 19

## Two Lecture Outline

**A quick introduction to…**

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

- Indexes
  - Inverted lists
  - Term dictionary
- Query processing
  - TAAT
  - DAAT
- Query operators

**Goal:** Provide an overview of search ("the Big Picture")
- Later lectures explore these topics in greater detail

39

© 2017, Jamie Callan

## Waitlist Reminder

**If you are on the waitlist…**
- I will admit some people from the waitlist today or tomorrow
- Be sure to sign the attendance sheet to show that you are here

157

© 2017, Jamie Callan