

---

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

**Introduction to Search**

Jamie Callan  
Carnegie Mellon University  
[callan@cs.cmu.edu](mailto: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
- 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

## Term Dictionary



### Main purpose

- Map term string to term id (integer)
- Example: “stocks” → 14,319

### Example entry for “stock”

id:	14,319
df:	1,070
idf:	1.146707
ctf:	3,567
inverted list:	<object handle>

### Often term dictionaries are stored in RAM

- So, they are a good place to store other frequently-access information, too
- E.g., statistics such as df, idf, ctf
- E.g., an inverted list pointer

40

© 2017, Jamie Callan

## Term Dictionary: Storage

### Example entry for “stock”

id:	14,319
df:	1,070
idf:	1.146707
ctf:	3,567
inverted list:	<object handle>

### Hash tables

- **O(1) lookup**
  - Very fast
- **Capabilities**
  - Exact-match lookup

### B-Trees

- **O (log n) lookup**
  - Fast, but not O(1)
- **Capabilities**
  - Exact-match lookup
  - Range lookup
  - Prefix lookup

41

© 2017, Jamie Callan

## Term Dictionary: Storage

Example entry for "block"

id:	14,319
df:	1,070
idf:	1.146707
ctf:	3,567
inverted list:	<object handle>

**The gov2 web corpus (25 million pages) has 40 million terms**

- The term complete dictionary might require 1.6 GB
- Most search engines won't devote that much RAM to a term dictionary

**Suppose we want to use a tiered term dictionary**

- Frequent terms in RAM (e.g.,  $ctf \geq 1,000$ )
- Less frequent terms on disk (e.g.,  $ctf < 1,000$ )
- The Indri search engine uses this strategy

**What percentage of terms are stored in RAM?**

42

© 2017, Jamie Callan

## Term Dictionary: Storage

Example entry for "block"

id:	14,319
df:	1,070
idf:	1.146707
ctf:	3,567
inverted list:	<object handle>

**What percentage of terms have  $ctf < 1,000$ ?**

- **Zipf's Law:**  $Rank_t \times Frequency_t = A \times N$
- **Rank of a word that occurs 1,000 times:**  $A \times N / 1000$
- **Rank of a word that occurs 1 time:**  $A \times N / 1$ 
  - An estimate of the vocabulary size
- **The percentage of terms with  $ctf < 1,000$ :**

$$(Rank_{Max} - Rank_{1000}) / Rank_{Max} =$$

$$((A \times N) - (A \times N / 1000)) / (A \times N) = 999 / 1000 = 99.9\%$$

**Ranking  
of terms  
by ctf**

**freq=1000** →

**freq=1** →



**But ... Zipf's Law is inaccurate at the tails**

- How does this affect the prediction?

43

© 2017, Jamie Callan

## Term Dictionary: Storage

Example entry for "blink"

id:	14,319
df:	1,070
idf:	1.146707
ctf:	3,567
inverted list:	<object handle>

### How good is the prediction on the gov2 corpus?

- 25 million .gov web pages with a vocabulary of 40 million terms

### How many terms occur more than 1,000 times?

**Predicted:** 40,000 frequent terms (0.10%)

**Actual:** 185,000 frequent terms (0.46%)

### About 99.5% of terms have $ctf \leq 1,000$

- 99.9% vs. 99.5%

44

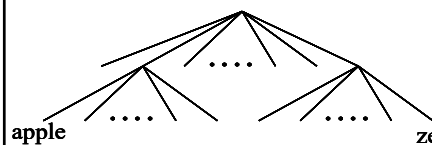
© 2017, Jamie Callan

## Term Dictionary: How Does Indri Do It?

Example entry for "blink"

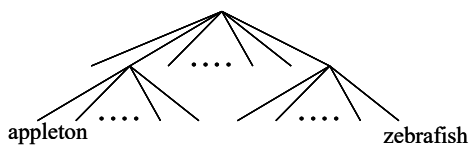
id:	14,319
df:	1,070
idf:	1.146707
ctf:	3,567
inverted list:	<object handle>

### Frequent Vocabulary ( $ctf \geq 1000$ )



**In Memory**

### Infrequent Vocabulary ( $ctf < 1000$ )



**Memory and Disk**

### For the Gov2 collection

- About 185,000 frequent terms in RAM (00.5%)
- About 39.78 million infrequent terms on disk (99.5%)

45

© 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
- 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

46

© 2017, Jamie Callan

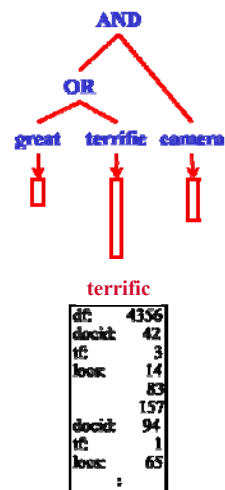
## Query Processing

### Query processing is a simple task

- **Input:** Two or more lists (e.g., inverted lists)
  - One list for each query argument
- **Output:** One list
  - The result list

### There are three approaches to query processing

- Term-at-a-Time (TAAT)
- Document-at-a-Time (DAAT)
- TAAT / DAAT hybrids



47

© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

### Key idea

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

48

© 2017, Jamie Callan

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

**Query:** (a AND b) OR c

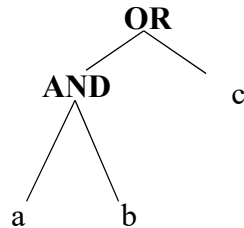
49

© 2017, Jamie Callan

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

**Query:** (a AND b) OR c

**Parsed  
Query**



50

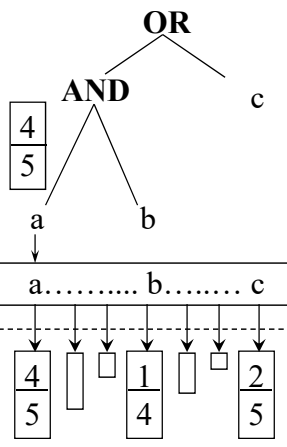
© 2017, Jamie Callan

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

**Query:** (a AND b) OR c

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

**Parsed  
Query**



**Term  
Dictionary**

**Inverted  
List  
Database**

51

© 2017, Jamie Callan

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

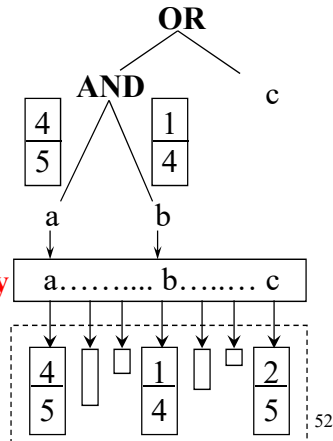
**Parsed  
Query**

**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  
Database**



© 2017, Jamie Callan

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

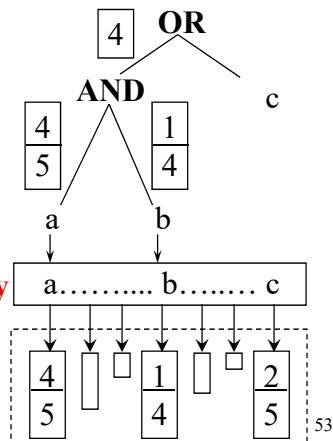
**Parsed  
Query**

**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
3. AND operator: Intersect the  
inverted lists for 'a' and 'b'

**Term  
Dictionary**

**Inverted  
List  
Database**



© 2017, Jamie Callan

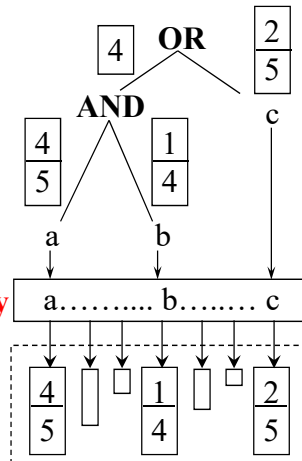


## Document Retrieval: Term-at-a-Time (TAAT) Query Evaluation

**Parsed  
Query**

**Term  
Dictionary**

**Inverted  
List  
Database**



**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

54

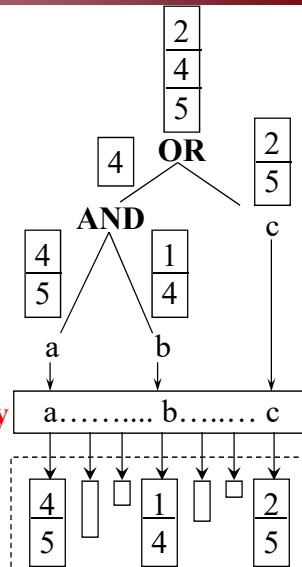
© 2017, Jamie Callan

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

**Parsed  
Query**

**Term  
Dictionary**

**Inverted  
List  
Database**



**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
5. OR operator: Union of AND operator results and 'c' inverted list

55

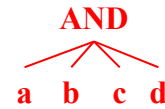
© 2017, Jamie Callan

## Query Processing: Term-at-a-Time 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
- $\text{Result}_{\text{AND\_2}} \text{ AND } d \rightarrow \text{Result}_Q$



56

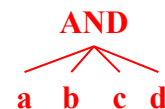
© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

**Query:** #AND (a b c d)

### Characteristics

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



57

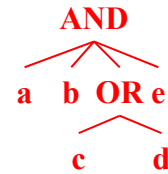
© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

**Query:** #AND (a b #OR (c d) e)

### Evaluation strategy

- Retrieve a
- Retrieve b
- $a \text{ AND } b \rightarrow \text{Result}_{\text{AND\_1}}$
- Retrieve c
- Retrieve d
- $c \text{ OR } d \rightarrow \text{Result}_{\text{OR}}$
- $\text{Result}_{\text{AND\_1}} \text{ AND } \text{Result}_{\text{OR}} \rightarrow \text{Result}_{\text{AND\_2}}$
- Retrieve e
- $\text{Result}_{\text{AND\_2}} \text{ AND } e \rightarrow \text{Result}_Q$



58

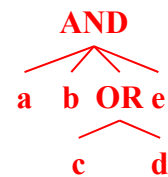
© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

**Query:** #AND (a b #OR (c d) e)

### Peak memory usage (probably)

- 4 lists in memory simultaneously
- $\text{size}(a \text{ AND } b) + \text{size}(c) + \text{size}(d) + \text{size}(c \text{ OR } d)$  bytes



59

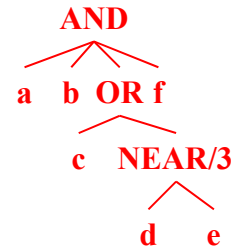
© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

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

### Evaluation strategy

- Retrieve a
- Retrieve b
- $a \text{ AND } b \rightarrow \text{Result}_{\text{AND\_1}}$
- Retrieve c
- Retrieve d
- Retrieve e
- $d \text{ NEAR/3 } e \rightarrow \text{Result}_{\text{NEAR}}$
- $c \text{ OR } \text{Result}_{\text{NEAR}} \rightarrow \text{Result}_{\text{OR}}$
- $\text{Result}_{\text{AND\_1}} \text{ AND } \text{Result}_{\text{OR}} \rightarrow \text{Result}_{\text{AND\_2}}$
- Retrieve f
- $\text{Result}_{\text{AND\_2}} \text{ AND } f \rightarrow \text{Result}_Q$



60

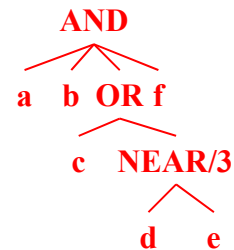
© 2017, Jamie Callan

## Query Processing: Term-at-a-Time 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



61

© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

### TAAT systems are easy to build

- Each operator is a series of  $f(\text{arg}_1, \text{arg}_2) \rightarrow \text{result}$  operations
- Conceptually, easy to understand
- Thus, we cover them first

### TAAT systems are very efficient

- Little wasted effort
  - This isn't apparent now, but will be when we consider DAAT

62

© 2017, Jamie Callan

## Query Processing: Term-at-a-Time Query Evaluation

### TAAT memory usage is uncontrolled

- Each query operator stores in RAM up to 3 lists simultaneously
  - $\text{arg}_1$ ,  $\text{arg}_2$ , result
- A query of depth  $d$  must store  $d+2$  lists in RAM simultaneously

### TAAT systems can run out of memory

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

### Rarely used in large-scale systems

63

© 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
- 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

64

© 2017, Jamie Callan

## Query Processing: Document-at-a-Time Query Evaluation

### Key idea

- Compute a complete score for  $\text{doc}_i$  before proceeding to  $\text{doc}_{i+1}$

**The following example assumes an unranked Boolean model.**

- All scores are 1
- The same architecture can be used for ranked Boolean

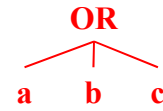
65

© 2017, Jamie Callan

## Query Processing: Document-at-a-Time Query Evaluation

### Starting condition:

- The query is #OR (a b c)
- Retrieve the inverted list for each term



<b>a:</b>	<b>b:</b>	<b>c:</b>
docid 19, score 1	docid 16, score 1	docid 17, score 1
docid 32, score 1	docid 19, score 1	docid 19, score 1
docid 42, score 1	docid 44, score 1	docid 44, score 1
docid 53, score 1	docid 51, score 1	docid 49, score 1
: :	: :	: :

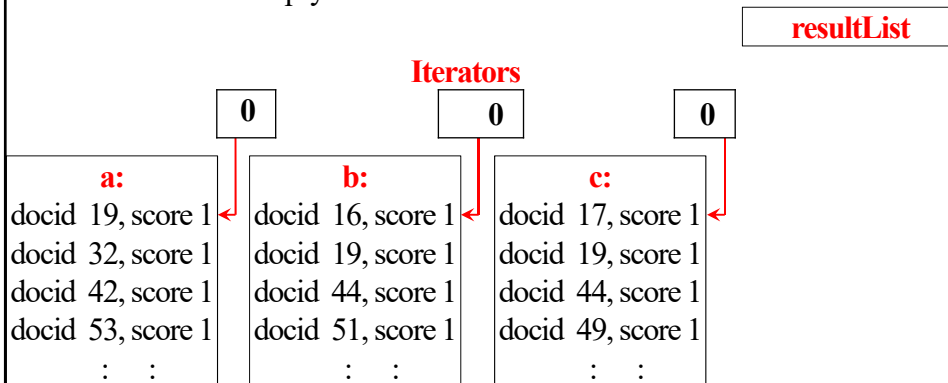
66

© 2017, Jamie Callan

## Query Processing: Document-at-a-Time Query Evaluation

### Initialization:

- Allocate iterators for processing inverted lists
- Allocate an empty result list



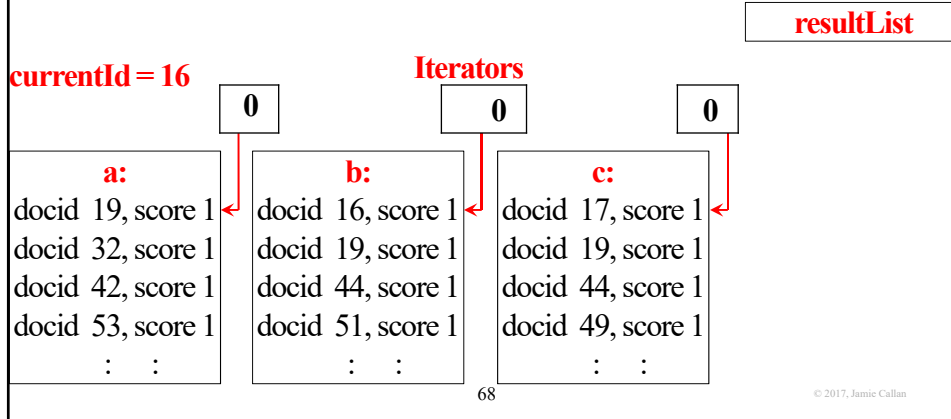
67

© 2017, Jamie Callan

## Query Processing: Document-at-a-Time Query Evaluation

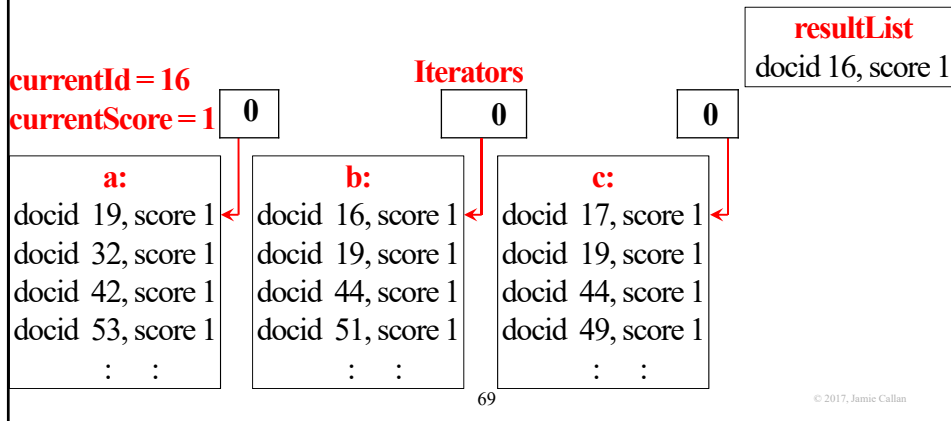
### Beginning of loop

- Examine the current document id in each list
- Find the minimum id



## Query Processing: Document-at-a-Time Query Evaluation

- Examine each list that contains currentId to compute currentScore
- Store the result

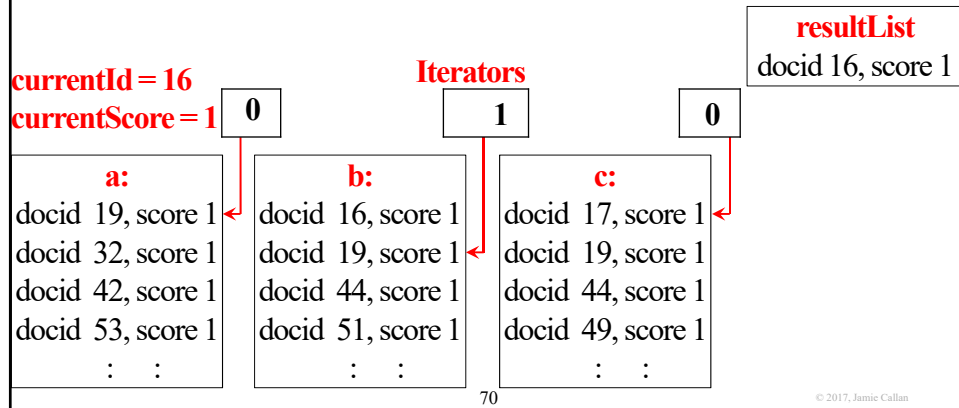




## Query Processing: Document-at-a-Time Query Evaluation

- Advance each pointer that points to the currentId

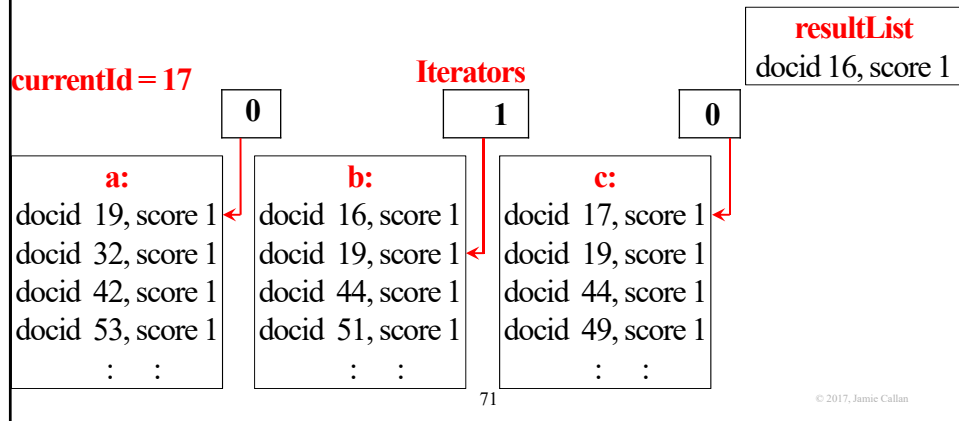
End of loop



## Query Processing: Document-at-a-Time Query Evaluation

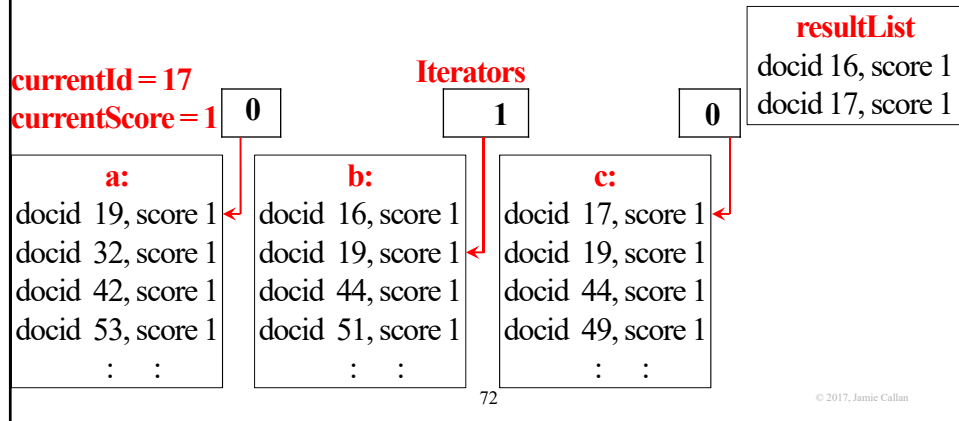
Beginning of loop

- Examine the current document id in each list
- Find the minimum id



## Query Processing: Document-at-a-Time Query Evaluation

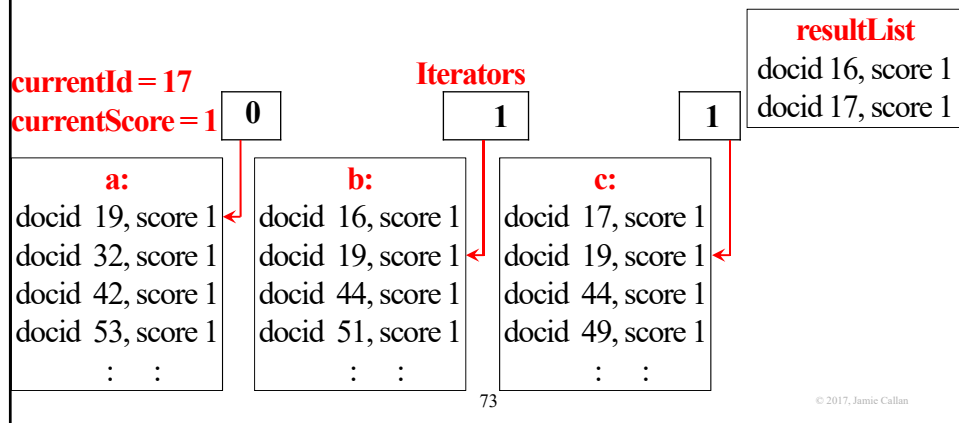
- Examine each list that contains currentId to compute currentScore
- Store the result



## Query Processing: Document-at-a-Time Query Evaluation

- Advance each pointer that points to the currentId

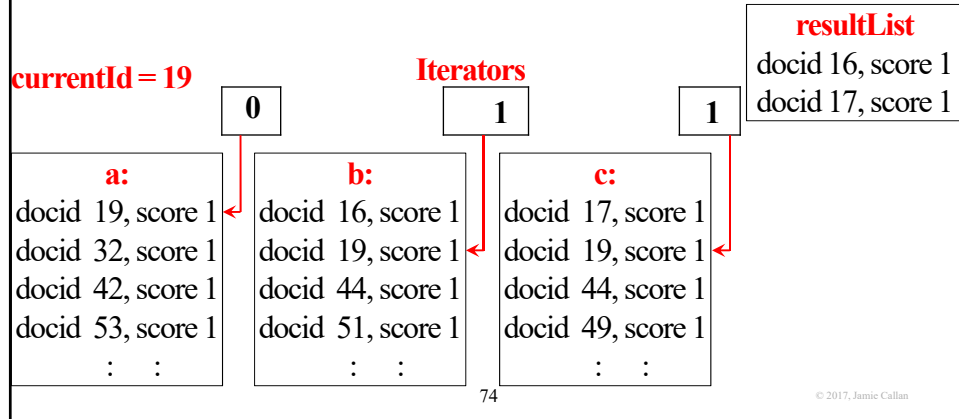
**End of loop**



## Query Processing: Document-at-a-Time Query Evaluation

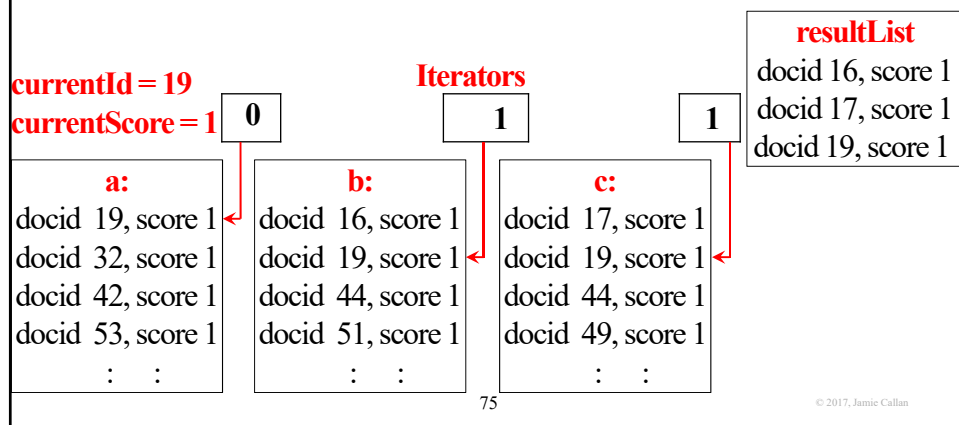
### Beginning of loop

- Examine the current document id in each list
- Find the minimum id



## Query Processing: Document-at-a-Time Query Evaluation

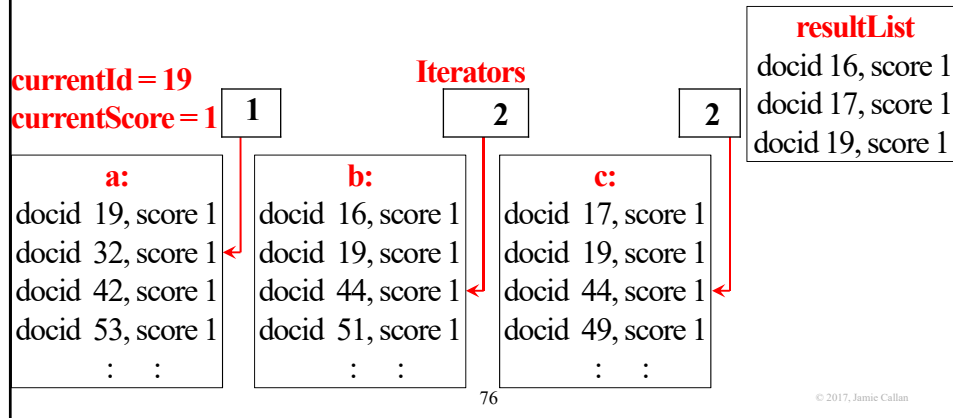
- Examine each list that contains currentId to compute currentScore
- Store the result



## Query Processing: Document-at-a-Time Query Evaluation

- Advance each pointer that points to the currentId

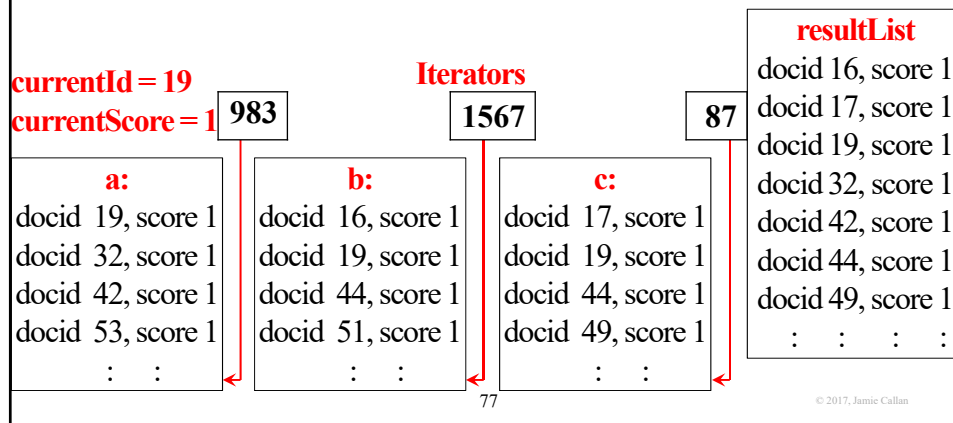
End of loop



## Query Processing: Document-at-a-Time Query Evaluation

Continue the loop until every inverted list is fully processed

Return the resultList



## DAAT Query Processing: Implementation Notes

### The simple implementation requires nested loops

- E.g., to find the minimum document id
- E.g., to compute the score for the current document id
- E.g., to advance pointers

### A more efficient implementation combines loops

- If this list has the current docid
  - Update the score
  - Advance the pointer

### There are many opportunities for clever optimization

78

© 2017, Jamie Callan

## Query Processing: Document-at-a-Time Query Evaluation

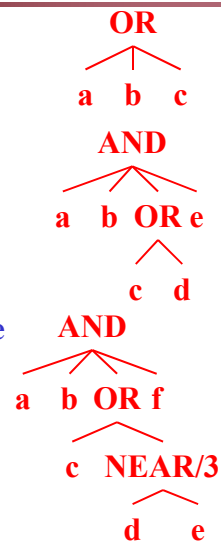
### How does DAAT support structured queries?

- Conceptually, it is something like this

```
q.initialize ()
while ( q.hasNext () )
    q.evalNext () returns next [docid, score] tuple
```

### Each call to `q.evalNext()` traverses the entire tree

- This is a little inefficient ... but not horrible
- The tricky part is figuring out the next docid
- Many opportunities for optimization



79

© 2017, Jamie Callan

## DAAT Query Processing Characteristics

### DAAT memory usage is easy to control

- It needs simultaneous access to all inverted lists (which seems bad)
- But ... inverted lists are read from disk into RAM in blocks
  - E.g., read the inverted list in 64MB blocks
- When the end of the current block is reached, read the next block
- The block size determines how much RAM the query uses

### Many query evaluation optimizations are possible

- E.g., only partial evaluation of documents with low scores

### Frequently used in large-scale systems

80

© 2017, Jamie Callan

## Query Processing: TAAT / DAAT Hybrids

### Hybrid TAAT and DAAT architectures are common

- To get a blend of efficiency and memory control
- E.g., block-based TAAT
  - Compute TAAT over blocks of document ids
- A popular research topic

81

© 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
- 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

82

© 2017, Jamie Callan

## Query Operators

### Usually exact match systems have rich query languages

- Ranking is weak or missing
- Query languages provide control over what is matched

### Today’s focus

- Proximity operators
  - NEAR/n and WINDOW/n
- Types (classes) of query operators

**The goal is to prepare you for HW1**

83

© 2017, Jamie Callan

## Proximity Operators: The NEAR Operator

**NEAR/n:** Distance between adjacent arguments is  $\geq 0$  &  $\leq n$  terms

- Query: “President NEAR/2 Obama”

Document texts:

“President Obama”	<b>Matches (distance is 1)</b>
“President Barack Obama”	<b>Matches (distance is 2)</b>
“President Barack H. Obama”	<b>Doesn’t match (distance is 3)</b>
“Obama is President”	<b>Doesn’t match (distance is -2)</b>

**Sentence/n:** Like NEAR/n, but distance is measured in sentences

**Paragraph/n:** Like NEAR/n, but distance is measured in paragraphs

84

© 2017, Jamie Callan

## Proximity Operators: The NEAR Operator

**The NEAR/n operator is used to match names and phrases**

- Arguments must be matched in order
- n specifies the maximum distance between adjacent terms

### Examples

- #NEAR/1 (barack obama)
  - Matches “barack obama”
  - Doesn’t match “barack hussein obama” or “obama, barack”
- #NEAR/3 (barack obama)
  - Matches “barack obama” and “barack hussein obama”
- #NEAR/4 (a b c) matches (a x x b x x x c)

85

© 2017, Jamie Callan



## Proximity Operators: A Simple Greedy NEAR/n Algorithm

Typically proximity operators have complexity  $O(|C|)$

- A single pass down each inverted list
- Similar in complexity to merging sorted lists

86

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

		Query: #NEAR/3 (a b)
a	b	
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

87

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b
df: 47	df: 95
<span style="border: 1px solid red; padding: 2px;">doc: 19</span>	<span style="border: 1px solid red; padding: 2px;">doc: 23</span>
tf: 1	tf: 1
locs: 7	locs: 99
doc: 27	doc: 27
tf: 3	tf: 4
locs: 47	locs: 48
98	49
132	133
doc: 92	134
...	doc: 148
	...

Query: #NEAR/3 (a b)

Initialize doc iterators

88

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b
df: 47	df: 95
<span style="border: 1px solid red; padding: 2px;">doc: 19</span>	<span style="border: 1px solid red; padding: 2px;">doc: 23</span>
tf: 1	tf: 1
locs: 7	locs: 99
doc: 27	doc: 27
tf: 3	tf: 4
locs: 47	locs: 48
98	49
132	133
doc: 92	134
...	doc: 148
	...

Query: #NEAR/3 (a b)

19 ≠ 23 (different doc)

Advance ptr of smallest  
doc

89

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	<span style="border: 1px solid red;">doc: 23</span>	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red;">doc: 27</span>	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

90

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	<span style="border: 1px solid red;">doc: 23</span>	<b>27 ≠ 23 (different doc)</b>
tf: 1	tf: 1	<b>Advance ptr of smallest</b>
locs: 7	locs: 99	<b>doc</b>
<span style="border: 1px solid red;">doc: 27</span>	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

91

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red;">doc: 27</span>	<span style="border: 1px solid red;">doc: 27</span>	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

92

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red;">doc: 27</span>	<span style="border: 1px solid red;">doc: 27</span>	<b>Same document</b>
tf: 3	tf: 4	<b>Initialize location ptrs</b>
<span style="border: 1px dashed blue;">locs: 47</span>	<span style="border: 1px dashed blue;">locs: 48</span>	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

93

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<b>48 – 47 ≤ n (match)</b>
tf: 3	tf: 4	<b>Record match</b>
<span style="border: 1px dashed blue; padding: 2px;">locs: 47</span>	<span style="border: 1px dashed blue; padding: 2px;">locs: 48</span>	<b>Advance all loc ptrs</b>
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

94

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	
tf: 3	tf: 4	
locs: 47	locs: 48	
<span style="border: 1px dashed blue; padding: 2px;">98</span>	<span style="border: 1px dashed blue; padding: 2px;">49</span>	
132	133	
doc: 92	134	
...	doc: 148	
	...	

95

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<b>98 &gt; 49 (no match)</b>
tf: 3	tf: 4	<b>Advance ptr of smallest</b>
locs: 47	locs: 48	<b>loc</b>
<span style="border: 1px dashed blue; padding: 2px;">98</span>	<span style="border: 1px dashed blue; padding: 2px;">49</span>	
132	133	
doc: 92	134	
...	doc: 148	
	...	

96

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	
tf: 3	tf: 4	
locs: 47	locs: 48	
<span style="border: 1px dashed blue; padding: 2px;">98</span>	49	
132	<span style="border: 1px dashed blue; padding: 2px;">133</span>	
doc: 92	134	
...	doc: 148	
	...	

97

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<b>133 – 98 &gt; n (no match)</b>
tf: 3	tf: 4	<b>Advance ptr of smallest</b>
locs: 47	locs: 48	<b>loc</b>
<span style="border: 1px dashed blue; padding: 2px;">98</span>	49	
132	<span style="border: 1px dashed blue; padding: 2px;">133</span>	
doc: 92	134	
...	doc: 148	
	...	

98

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>Query: #NEAR/3 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
<span style="border: 1px dashed blue; padding: 2px;">132</span>	<span style="border: 1px dashed blue; padding: 2px;">133</span>	
doc: 92	134	
...	doc: 148	
	...	

99

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>
df: 47	df: 95
doc: 19	doc: 23
tf: 1	tf: 1
locs: 7	locs: 99
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>
tf: 3	tf: 4
locs: 47	locs: 48
98	49
<span style="border: 1px dashed blue; padding: 2px;">132</span>	<span style="border: 1px dashed blue; padding: 2px;">133</span>
doc: 92	134
...	doc: 148
	...

**Query: #NEAR/3 (a b)**

**$133 - 132 \leq n$  (match)**

**Record match**

**Advance all loc ptrs**

**a's locs are exhausted.**

**No more matches are  
possible in this document**

**Advance all doc ptrs.**

100

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>
df: 47	df: 95
doc: 19	doc: 23
tf: 1	tf: 1
locs: 7	locs: 99
doc: 27	doc: 27
tf: 3	tf: 4
locs: 47	locs: 48
98	49
132	133
<span style="border: 1px solid red; padding: 2px;">doc: 92</span>	134
...	<span style="border: 1px solid red; padding: 2px;">doc: 148</span>
	...

**Query: #NEAR/3 (a b)**

**Continue until the inverted  
lists are exhausted**

101

© 2017, Jamie Callan



## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>c</b>	<b>Query: #NEAR/3 (a b c)</b>
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

102

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

<b>a</b>	<b>b</b>	<b>c</b>	<b>Query: #NEAR/3 (a b c)</b>
df: 47	df: 95	df: 14	
<b>doc: 19</b>	<b>doc: 23</b>	<b>doc: 23</b>	<b>Initialize doc iterators</b>
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

103

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
<span style="border: 1px solid red; padding: 2px;">doc: 19</span>	<span style="border: 1px solid red; padding: 2px;">doc: 23</span>	<span style="border: 1px solid red; padding: 2px;">doc: 23</span>	<b>19 ≠ 23 (different doc)</b>
tf: 1	tf: 1	tf: 1	<b>Advance ptr of smallest doc</b>
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

104

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	<span style="border: 1px solid red; padding: 2px;">doc: 23</span>	<span style="border: 1px solid red; padding: 2px;">doc: 23</span>	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

105

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	<b>27 ≠ 23 (different doc)</b>
tf: 1	tf: 1	tf: 1	<b>Advance ptrs of smallest docs</b>
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

106

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

107

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	Same document
tf: 1	tf: 1	tf: 1	Initialize location ptrs
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

108

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	$48 - 47 \leq n$ (match)
tf: 3	tf: 4	tf: 4	$46 - 48 < 0$ (no match)
locs: 47	locs: 48	locs: 46	Advance ptr of smallest loc
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

109

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

110

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

111

© 2017, Jamie Callan

**$48 - 47 \leq n$  (match)**  
 **$51 - 48 \leq n$  (match)**  
**Record match**  
**Advance all loc ptrs**

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

112

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

$49 - 98 < 0$  (no match)  
 Advance ptr of smallest loc

113

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

114

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

115

© 2017, Jamie Callan

**133 – 98 > n (no match)**  
**Advance ptr of smallest loc**

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

116

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

**$133 - 132 > n$  (match)**  
 **$114 - 133 < 0$  (no match)**  
**Advance ptr of smallest loc**

117

© 2017, Jamie Callan



## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

118

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	
doc: 19	doc: 23	doc: 23	
tf: 1	tf: 1	tf: 1	
locs: 7	locs: 99	loc: 99	
doc: 27	doc: 27	doc: 27	
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	...	...	

$133 - 132 \leq n$  (match)  
 $137 - 133 > n$  (no match)  
 Advance ptr of smallest loc  
 a's locs are exhausted.  
 No more matches are possible in this document  
 Advance all doc ptrs

119

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c
df: 47	df: 95	df: 14
doc: 19	doc: 23	doc: 23
tf: 1	tf: 1	tf: 1
locs: 7	locs: 99	loc: 99
doc: 27	doc: 27	doc: 27
tf: 3	tf: 4	tf: 4
locs: 47	locs: 48	locs: 46
98	49	51
132	133	114
doc: 92	134	137
...	doc: 148	doc: 129

Query: #NEAR/3 (a b c)

Perhaps you expected b's  
loc ptr to be advanced

This is a flaw in the simple  
greedy algorithm. The  
smallest loc ptr isn't always  
the optimal loc ptr

Better algorithms are  
possible, but also more  
complex

120

© 2017, Jamie Callan

## Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c
df: 47	df: 95	df: 14
doc: 19	doc: 23	doc: 23
tf: 1	tf: 1	tf: 1
locs: 7	locs: 99	loc: 99
doc: 27	doc: 27	doc: 27
tf: 3	tf: 4	tf: 4
locs: 47	locs: 48	locs: 46
98	49	51
132	133	114
doc: 92	134	137
...	doc: 148	doc: 129

Query: #NEAR/3 (a b c)

Continue until the inverted  
lists are exhausted

121

© 2017, Jamie Callan

## Proximity Operators: NEAR/n FAQ

**Query:** #NEAR/2 (a b)

**Text:** a b x a x x x a x x b x x a x b

- There are two matches {1, 2} and {14, 16}
- Results for the NEAR operator: tf=2, and locations=2, 16

**Query:** #NEAR/2 (a b c)

**Text:** a a b b c c

- There are two matches {1, 3, 5} and {2, 4, 6}
- Results for the NEAR operator: tf=2, and locations=5, 6

122

© 2017, Jamie Callan

## Proximity Operators: NEAR/n FAQ

**Query:** #NEAR/3 (a b)

**Text:** a b c b

- There is one match {1, 2}
  - A query term can match only one text term
- Results for the NEAR operator: tf=1, and locations=2

**Query:** #NEAR/3 (a b)

**Text:** b a c a

- There are no matches
  - The order of NEAR query arguments is important

123

© 2017, Jamie Callan

## Proximity Operators: NEAR/n FAQ

**Query:** #NEAR/2 (a a b)

**Text:** a b a a b

- The greedy algorithm does not handle duplicate arguments
- The distance constraint can be modified to support this case
  - E.g.,  $0 < \text{distance}(t_1, t_2) \leq n$
  - Maybe also a better algorithm for advancing pointers

**These cases haven't been studied much by researchers**

- But Google seems to support “apple apple pie”

124

© 2017, Jamie Callan

## Proximity Operators: NEAR/n FAQ

**Query:** #NEAR/3 (a b c)

**Text:** a b d b x x c

- The greedy algorithm fails to find a match
  - It considers {1, 2, 7} (the first location for each term)
  - {1, 2, 7} fails to match, so the smallest location pointer advances
    - » a
  - The list for a is exhausted, so this text does not match
- A better algorithm would find a match at {1, 4, 7}
  - More accurate, but also slower
  - The greedy algorithm is usually sufficient in practice
- Use the greedy algorithm for your homework

125

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

The **WINDOW/n** operator is used to match related concepts

- Arguments can be in any order
- n specifies the maximum distance between any pair of terms

### Examples

- WINDOW/100 (obama merkel putin)
  - We don't care which order they occur in

Typically proximity operators have complexity  $O(|C|)$

- A single pass down each inverted list

126

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

		Query: # WINDOW/20 (a b)
a	b	
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

127

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	Initialize doc iterators
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

128

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	19 ≠ 23 (different doc)
tf: 1	tf: 1	Advance ptr of smallest doc
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

129

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	<span style="border: 1px solid red;">doc: 23</span>	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red;">doc: 27</span>	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

130

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	<span style="border: 1px solid red;">doc: 23</span>	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red;">doc: 27</span>	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

**27 ≠ 23 (different doc)**  
**Advance ptr of smallest doc**

131

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

132

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	Same document
tf: 3	tf: 4	Initialize location ptrs
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

133

© 2017, Jamie Callan



## Proximity Operators: The Window (or Unordered Window) Operator

<b>a</b>	<b>b</b>	<b>Query: # WINDOW/20 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<b> 48 – 47  &lt; 20 (match)</b>
tf: 3	tf: 4	<b>Record match</b>
<span style="border: 1px dashed blue; padding: 2px;">locs: 47</span>	<span style="border: 1px dashed blue; padding: 2px;">locs: 48</span>	<b>Advance all loc ptrs</b>
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

134

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

<b>a</b>	<b>b</b>	<b>Query: # WINDOW/20 (a b)</b>
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	<span style="border: 1px solid red; padding: 2px;">doc: 27</span>	
tf: 3	tf: 4	
locs: 47	locs: 48	
<span style="border: 1px dashed blue; padding: 2px;">98</span>	<span style="border: 1px dashed blue; padding: 2px;">49</span>	
132	133	
doc: 92	134	
...	doc: 148	
	...	

135

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	98 – 49  ≥ 20 (no match)
132	133	Advance ptr of smallest
doc: 92	134	loc
...	doc: 148	
	...	

136

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

137

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	133 – 98  ≥ 20 (no match)
tf: 3	tf: 4	Advance ptr of smallest
locs: 47	locs: 48	loc
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

138

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b	Query: # WINDOW/20 (a b)
df: 47	df: 95	
doc: 19	doc: 23	
tf: 1	tf: 1	
locs: 7	locs: 99	
doc: 27	doc: 27	
tf: 3	tf: 4	
locs: 47	locs: 48	
98	49	
132	133	
doc: 92	134	
...	doc: 148	
	...	

139

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b
df: 47	df: 95
doc: 19	doc: 23
tf: 1	tf: 1
locs: 7	locs: 99
doc: 27	doc: 27
tf: 3	tf: 4
locs: 47	locs: 48
98	49
132	133
doc: 92	134
...	doc: 148
	...

Query: # WINDOW/20 (a b)

$|133 - 132| < 20$  (match)

Record match

Advance all loc ptrs

a's locs are exhausted.

No more matches are  
possible in this document

Advance all doc ptrs.

140

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

a	b
df: 47	df: 95
doc: 19	doc: 23
tf: 1	tf: 1
locs: 7	locs: 99
doc: 27	doc: 27
tf: 3	tf: 4
locs: 47	locs: 48
98	49
132	133
doc: 92	134
...	doc: 148
	...

Query: # WINDOW/20 (a b)

Continue until the inverted  
lists are exhausted

141

© 2017, Jamie Callan

## Proximity Operators: The Window (or Unordered Window) Operator

### Implementation note

- A document term can only match the query once
- **Query:** #WINDOW/100 (obama merkel putin)
- **Document:** obama ... merkel ... putin ... merkel ... obama
- There is just one match here

### One can imagine other implementations, but this is the norm

- Usually more complicated matching doesn't improve accuracy

142

© 2017, Jamie Callan

## Two Important Data Structures

### Inverted lists

'apple'

df:	27
docid:	14
tf:	2
locs:	37
	92
docid:	89
	:
	:

**You already  
know about  
inverted lists  
(stored on disk)**

### Score lists

'apple'

<u>docid</u>	<u>score</u>
14	3
89	2
127	2
:	:
:	:

**#AND and #OR  
query operators  
produce these  
(in memory)**

**An unranked  
system has  
score = 1**

**How do we get from inverted lists to score lists?**

143

© 2017, Jamie Callan

## Three Types of Query Operators

There are three important classes of query operators

- Operators that combine inverted lists
  - E.g., #NEAR/n, #SYNONYM, #WINDOW/n, ...
  - Dynamically create new index terms / concepts
- Operators that convert inverted lists to scores
  - Usually invisible or implied
  - Usually just one – we will call it #SCORE
- Operators that combine score lists
  - E.g., #AND, #OR, #AVERAGE, #MAX, ...
  - Combine evidence

**invertedList**  
df  
docid tf loc+  
docid tf loc+  
: : :

**scoreList**  
df  
docid score  
docid score  
...

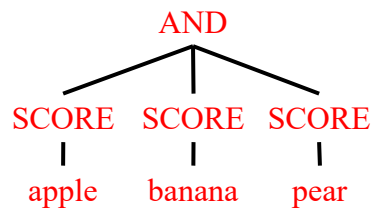
Every (?) query operator belongs to one of these classes

144

© 2017, Jamie Callan

## Three Types of Query Operators

**#AND (apple banana pear)**

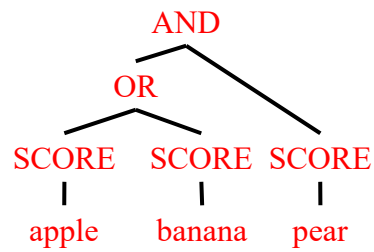


145

© 2017, Jamie Callan

## Three Types of Query Operators

**#AND (#OR (apple banana) pear)**

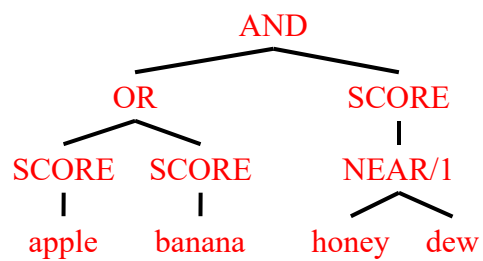


146

© 2017, Jamie Callan

## Three Types of Query Operators

**#AND (#OR (apple banana) #NEAR/1 (honey dew))**



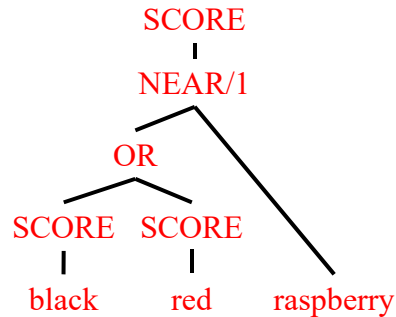
147

© 2017, Jamie Callan

## Three Types of Query Operators

Some queries are not allowed

**#SCORE (#NEAR/1 (#OR (black red) raspberry))**



**NEAR/1 operates on  
inverted lists**

It creates an inverted list  
for a new concept

**OR operates on score lists**  
It combines evidence

**Use a SYNONYM operator  
to combine inverted lists**

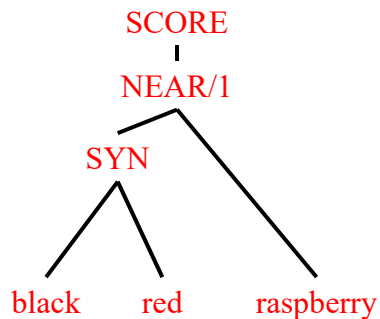
148

© 2017, Jamie Callan

## Three Types of Query Operators

This query is allowed

**#SCORE (#NEAR/1 (#SYN (black red) raspberry))**



**NEAR/1 operates on  
inverted lists**

It creates an inverted list  
for a new concept

149

© 2017, Jamie Callan



## OR vs SYN (SYNONYM)

OR and SYN look similar, but they are very different

SYN constructs new concepts dynamically

- By merging inverted lists

OR combines evidence about how well the document satisfies the information need

- Evidence obtained from matching multiple terms

These are very different things

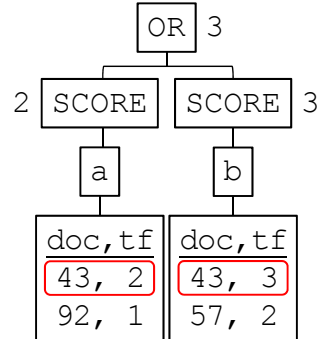
150

© 2017, Jamie Callan

## OR vs SYN (SYNONYM): Unranked Boolean Retrieval Model

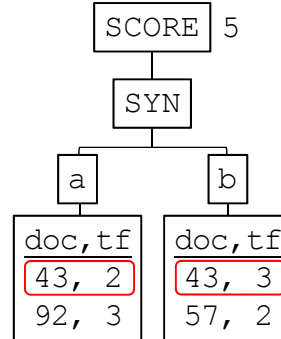
#OR (a b)

Matches two terms



#SYN (a b)

Matches one combined term



**SCORE:** tf

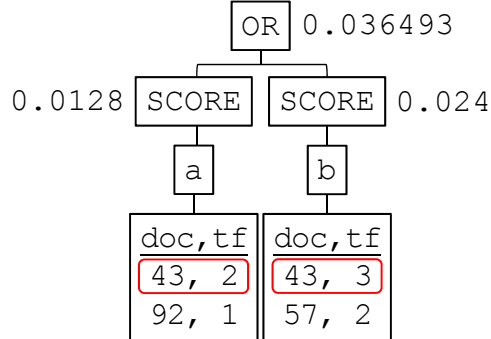
**OR:** MAX

© 2017, Jamie Callan

## OR vs SYN (SYNONYM): Indri Retrieval Model

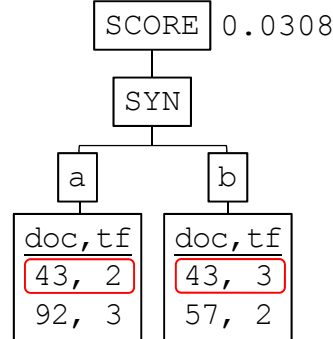
#OR (a b)

Matches two terms



#SYN (a b)

Matches one combined term



$$\text{SCORE: } p_{\text{score}}(q_i | d) = \frac{tf_{q_i, d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu}$$

$$\text{OR: } p_{or}(q | d) = 1 - \prod_{q_i \in q} (1 - p(q_i | d))$$

© 2017, Jamie Callan

## OR vs SYN (SYNONYM)

**SYN behaves the same for all retrieval models**

- i.e., merge inverted lists to create new concepts

**Retrieval models treat the concepts produced by SYN differently**

- **Unranked Boolean:** Prefers concepts that are frequent in this document
  - tf-oriented
- **Indri:** Prefers concepts that are frequent in this document, but rare across the corpus
  - tf.idf-oriented

153

© 2017, Jamie Callan

## Do Real Systems Have a SCORE Operator?

People don't write SCORE operators in their queries  
... so how do SCORE operators get into queries?

The query parser can insert them automatically when needed

- If a query operator combines scores &&  
its  $i^{\text{th}}$  query argument returns an inverted list  
Then wrap the  $i^{\text{th}}$  argument in a SCORE operator
- E.g., #AND (a b)  $\rightarrow$  #AND (#SCORE (a) #SCORE (b))
- The sample QryEval software does this

154

© 2017, Jamie Callan

## Outline

### 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
  - TAAT / DAAT hybrids
- Query operators

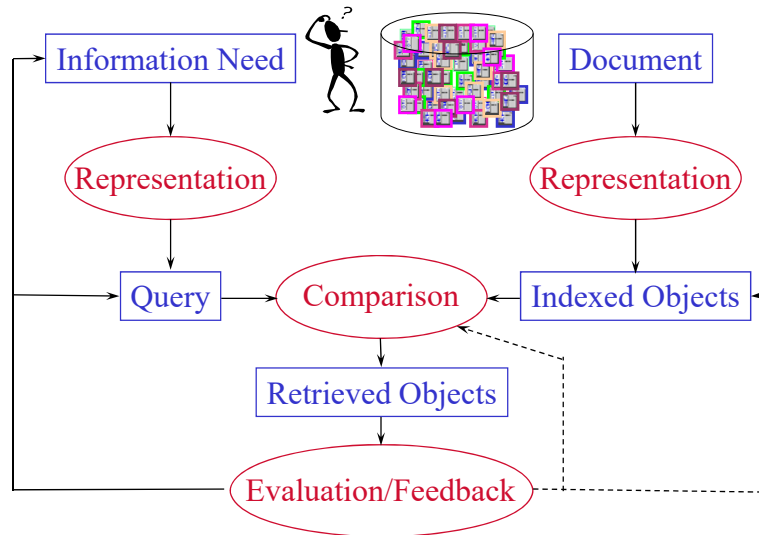
**Goal:** Provide an overview of search (“the Big Picture”)

- Later lectures explore these topics in greater detail

155

© 2017, Jamie Callan

## Overview of Information Retrieval Processes



156

© 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