
**11-442 / 11-642 / 11-742:
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

54

© 2021, Jamie Callan

54

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Key idea

- Compute a complete score for doc_i before proceeding to doc_{i+1}

The following example assumes an unranked Boolean model.

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

55

© 2021, Jamie Callan

55

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Starting condition:

- The query is #OR (a b c)

56

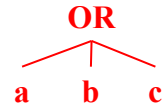
© 2021, Jamie Callan

56

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Starting condition:

- The query is #OR (a b c)
- Parse the query



57

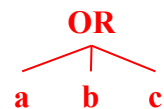
© 2021, Jamie Callan

57

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Starting condition:

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



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

58

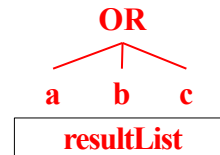
© 2021, Jamie Callan

58

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Initialization:

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



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

59

© 2021, Jamie Callan

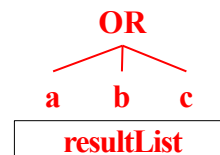
59

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Beginning of loop

- Examine the visible document id in each list
- Set currentId to the minimum id

currentId = 16



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

60

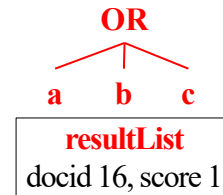
© 2021, Jamie Callan

60

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

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

currentId = 16
currentScore = 1



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

61

© 2021, Jamie Callan

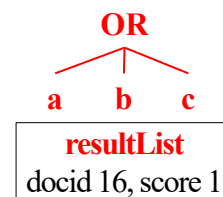
61

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Advance each iterator that points to the currentId

End of loop

currentId = 16
currentScore = 1



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

62

© 2021, Jamie Callan

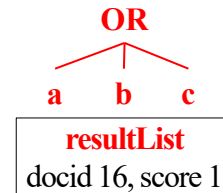
62

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Beginning of loop

- Examine the visible document id in each list
- Set currentId to the minimum id

currentId = 17



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

63

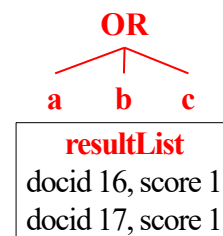
© 2021, Jamie Callan

63

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

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

currentId = 17
currentScore = 1



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

64

© 2021, Jamie Callan

64

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

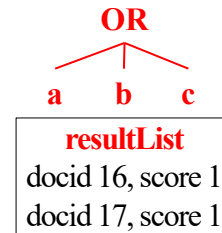
- Advance each iterator that points to the currentId

End of loop

currentId = 17
currentScore = 1

a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

65



© 2021, Jamie Callan

65

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

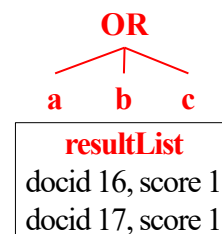
Beginning of loop

- Examine the visible document id in each list
- Set currentId to the minimum id

currentId = 19

a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

66



© 2021, Jamie Callan

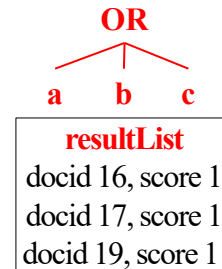
66

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

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

currentId = 19
currentScore = 1

a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:



67

© 2021, Jamie Callan

67

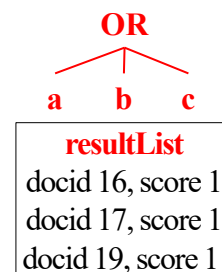
Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Advance each iterator that points to the currentId

End of loop

currentId = 19
currentScore = 1

a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:



68

© 2021, Jamie Callan

68

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Continue the loop until every inverted list is fully processed
Return the resultList

currentId = ...
currentScore = ...

a:
docid 19
docid 32
docid 42
docid 53
: :

b:
docid 16
docid 19
docid 44
docid 51
: :

c:
docid 17
docid 19
docid 44
docid 49
: :

OR			
a	b	c	
resultList			
docid 16, score 1			
docid 17, score 1			
docid 19, score 1			
docid 32, score 1			
docid 42, score 1			
docid 44, score 1			
docid 49, score 1			
:	:	:	:

69

© 2021, Jamie Callan

69

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Starting condition:

- The query is #AND (a b c)

70

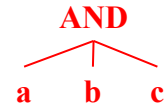
© 2021, Jamie Callan

70

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Starting condition:

- The query is #AND (a b c)
- Parse the query



71

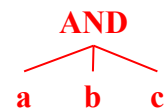
© 2021, Jamie Callan

71

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Starting condition:

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



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

72

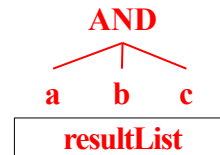
© 2021, Jamie Callan

72

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Initialization:

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



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

73

© 2021, Jamie Callan

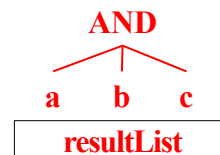
73

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Beginning of loop

- Examine the visible document id in each list
- Set currentId to the minimum id

currentId = 16



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

74

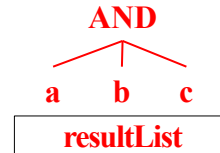
© 2021, Jamie Callan

74

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Examine each list that contains currentId to compute currentScore
 - No match

currentId = 16
currentScore = 0



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

75

© 2021, Jamie Callan

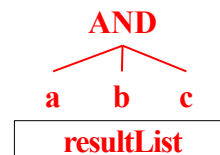
75

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Advance each iterator that points to the currentId

End of loop

currentId = 16
currentScore = 0



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

76

© 2021, Jamie Callan

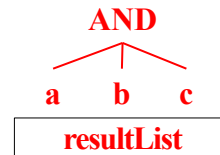
76

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Beginning of loop

- Examine the visible document id in each list
- Set currentId to the minimum id

currentId = 17



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

77

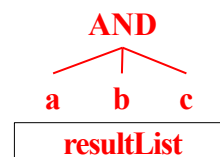
© 2021, Jamie Callan

77

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Examine each list that contains currentId to compute currentScore
 - No match

currentId = 17
currentScore = 0



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
:	:	:

78

© 2021, Jamie Callan

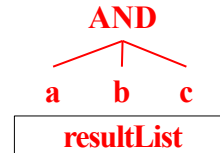
78

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Advance each iterator that points to the currentId

End of loop

currentId = 17
currentScore = 0



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

79

© 2021, Jamie Callan

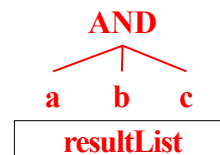
79

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Beginning of loop

- Examine the visible document id in each list
- Set currentId to the minimum id

currentId = 19



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

80

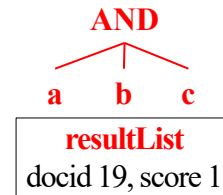
© 2021, Jamie Callan

80

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

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

currentId = 19
currentScore = 1



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

81

© 2021, Jamie Callan

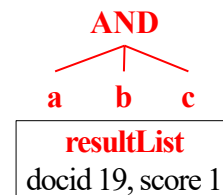
81

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

- Advance each iterator that points to the currentId

End of loop

currentId = 19
currentScore = 1



a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

82

© 2021, Jamie Callan

82

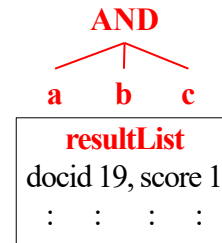
Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

Continue the loop until every inverted list is fully processed
Return the resultList

currentId = ...
currentScore = ...

a:	b:	c:
docid 19	docid 16	docid 17
docid 32	docid 19	docid 19
docid 42	docid 44	docid 44
docid 53	docid 51	docid 49
: :	: :	: :

83



© 2021, Jamie Callan

83

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

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 decide which iterators to advance

A more efficient implementation combines loops

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

There are many opportunities for clever optimization

84

© 2021, Jamie Callan

84

Document Retrieval: Document-at-a-Time (DAAT) Query Evaluation

How does DAAT support structured queries?

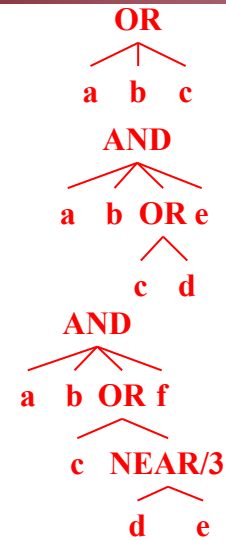
- Conceptually, it is something like this

```
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

The next lecture covers this in more detail



© 2021, Jamie Callan

85

85

DAAT Query Evaluation 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 256MB 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

86

© 2021, Jamie Callan

86

Document Retrieval: 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

87

© 2021, Jamie Callan

87

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

88

© 2021, Jamie Callan

88

Query Operators

Usually search engines have rich query languages

- Query languages provide control over what is matched

Today's focus

- Types (classes) of query operators
 - There are many operators, but just a few types of operators
- The NEAR/n proximity operator

The goal is to prepare you for HW1

89

© 2021, Jamie Callan

89

Three Types of Query Operators

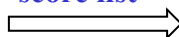
1. Produce new inverted lists

- Dynamically create new index terms / concepts
- E.g., #SYNONYM, #NEAR

Inverted list

df:	4356
docid:	42
tf:	3
locs:	14
	93
	157
:	

2. Use an inverted list to produce a score list



3. Combine scores

- Combine estimates about how well a document matches
- E.g., #AND, #OR, WSUM

Score list

length:	4356
docid: 42,	score: 3
docid: 89,	score: 5
docid: 127,	score: 4
:	:

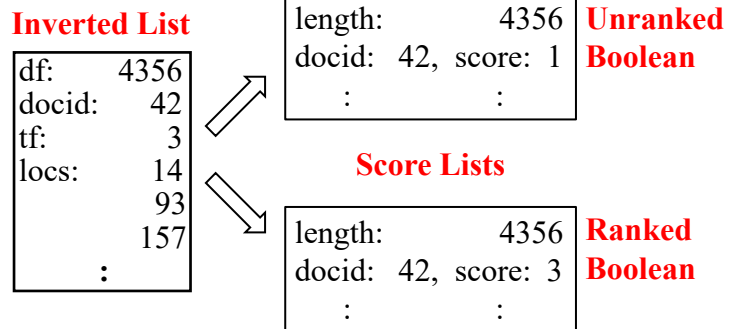
90

© 2021, Jamie Callan

90

Three Types of Query Operators: SCORE operator

The behavior of the SCORE operator depends on the retrieval model



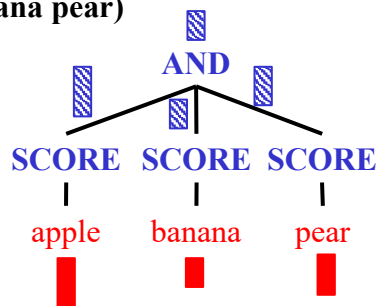
91

© 2021, Jamie Callan

91

Three Types of Query Operators

#AND (apple banana pear)



92

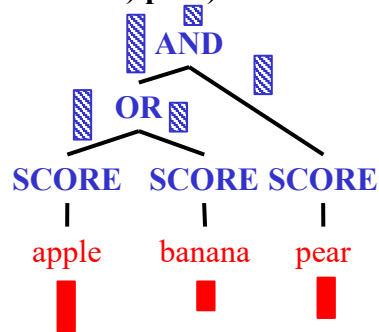
 **Inverted List**  **Score List**

© 2021, Jamie Callan

92

Three Types of Query Operators

#AND (#OR (apple banana) pear)



 Inverted List  Score List

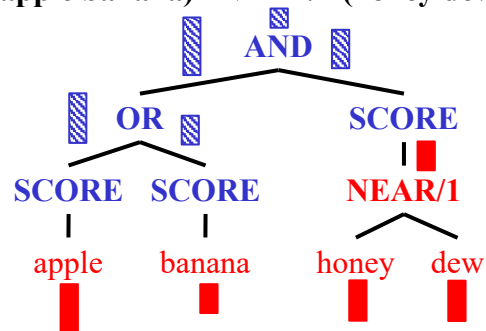
93

© 2021, Jamie Callan

93

Three Types of Query Operators

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



 Inverted List  Score List

94

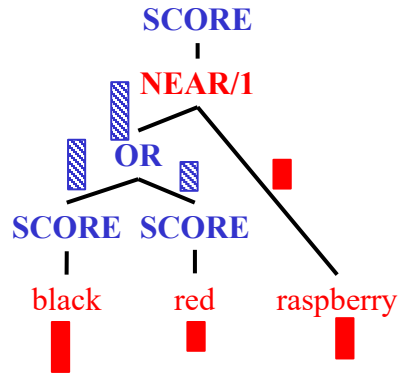
© 2021, Jamie Callan

94

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 (scores)
It does not produce locations

 Inverted List  Score List

95

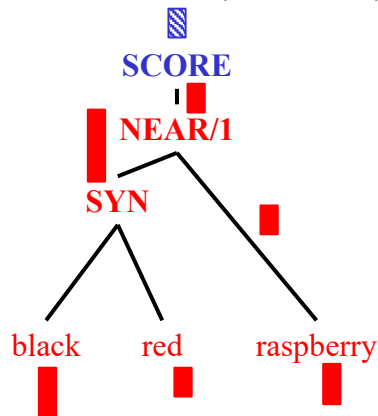
© 2021, Jamie Callan

95

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

The SYNONYM operator
combines inverted lists

It creates an inverted list
for a new concept

 Inverted List  Score List

96

© 2021, Jamie Callan

96

OR vs SYN (SYNONYM)

OR and SYN are very different

SYN dynamically constructs new concepts (new inverted lists)

- E.g., #SYN (red yellow green blue) represents the concept “color”
- The result is an inverted list showing where the concept occurs

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

- Each term's score is evidence about how well q matches d
- The result is a score list showing how well each document matches q

These operators produce different search engine behavior

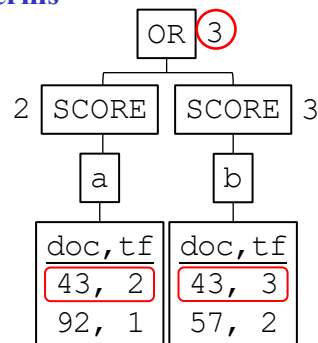
97

© 2021, Jamie Callan

97

OR vs SYN (SYNONYM): Ranked Boolean Retrieval Model

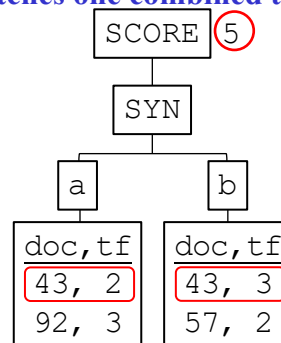
#OR (a b)
Matches two terms



SCORE: tf

OR: MAX

#SYN (a b)
Matches one combined term



#SYN (a b) > #OR (a b)

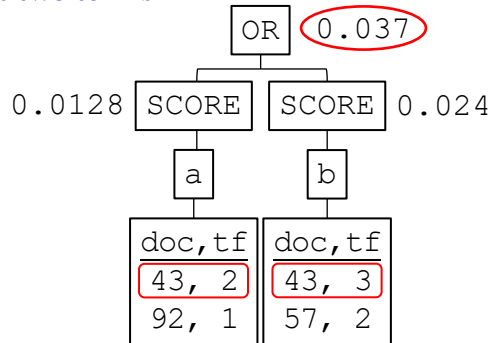
98

© 2021, Jamie Callan

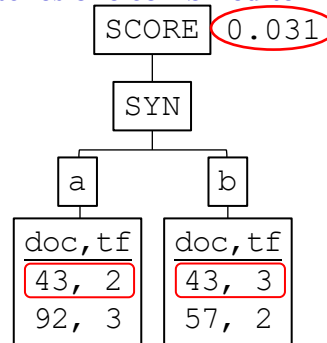
98

OR vs SYN (SYNONYM): Indri Retrieval Model

#OR (a b)
Matches two terms



#SYN (a b)
Matches one combined term



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

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

#OR (a b) > #SYN (a b)

© 2021, Jamie Callan

99

Where Do SCORE Operators Come From?

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

The query parser inserts them automatically

- If the query operator expects a score list argument && its argument is an operator that produces inverted lists
Then wrap the argument in a SCORE operator
- E.g., #AND (a b) → #AND (#SCORE (a) #SCORE (b))
- The QryEval homework software does this

100

© 2021, Jamie Callan

100

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

101

© 2021, Jamie Callan

101

Proximity Operators: The NEAR Operator

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

- Query: “President NEAR/2 Obama”

Document texts:

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

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

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

102

© 2021, Jamie Callan

102

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)

103

© 2021, Jamie Callan

103

Proximity Operators: A Simple Greedy NEAR/n Algorithm

There are many ways to implement the NEAR/n operator

- An exact implementation has high computational complexity
- Most implementations are greedy and inexact

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

- A single pass down each inverted list
- Similar in complexity to merging sorted lists
- They may not find some matches, but good enough for most tasks

Your implementation must match Jamie’s greedy algorithm

104

© 2021, Jamie Callan

104

Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	Query: #NEAR/3 (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	
	...	

105

© 2021, Jamie Callan

105

Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	Query: #NEAR/3 (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	
	...	

106

© 2021, Jamie Callan

106

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

**Advance all doc iterators
until they point to the
same document**

- This is a simple nested loop

107

© 2021, Jamie Callan

107

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

**Same document
Initialize location iterators**

108

© 2021, Jamie Callan

108

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

**Advance loc iterators $q_{i>0}$
left-to-right such that
 $\text{loc}(q_i) < \text{loc}(q_{i+1})$**
• Not necessary here

109

© 2021, Jamie Callan

109

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

**Do a left-to-right check
of locations**

$48 - 47 \leq n$ (match)

Record match

• **Right-most matching loc
(48)**

110

© 2021, Jamie Callan

110

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

Increment all loc iterators

111

© 2021, Jamie Callan

111

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

Advance loc iterators $q_{i>0}$
left-to-right such that
 $\text{loc}(q_i) < \text{loc}(q_{i+1})$

112

© 2021, Jamie Callan

112

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

Do a left-to-right check
of locations

$133 - 98 > 3$ (no match)

113

© 2021, Jamie Callan

113

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

Increment q_0 loc iterator

114

© 2021, Jamie Callan

114

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

**Advance loc iterators $q_{i>0}$
left-to-right such that
 $\text{loc}(q_i) < \text{loc}(q_{i+1})$**
• Not necessary here

115

© 2021, Jamie Callan

115

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

**Do a left-to-right check
of locations**

$133 - 132 \leq n$ (match)

Record match

• **Right-most matching loc
(133)**

116

© 2021, Jamie Callan

116

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

Increment all loc iterators

q_0 locs are exhausted.

No more matches are possible in this document

117

© 2021, Jamie Callan

117

Proximity Operators: A Simple Greedy NEAR/n Algorithm

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: #NEAR/3 (a b)

Increment all doc iterators

...

Continue until the inverted lists are exhausted

118

© 2021, Jamie Callan

118

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	
	
		119	

© 2021, Jamie Callan

119

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	Initialize doc iterators
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	
	
		120	

© 2021, Jamie Callan

120

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)

Advance all doc iterators
until they point to the
same document

- This is a simple nested loop

121

© 2021, Jamie Callan

121

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)

Same document
Initialize location iterators

122

© 2021, Jamie Callan

122

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	Advance loc iterators $q_{i>0}$
tf: 3	tf: 4	tf: 4	left-to-right such that
locs: 47	locs: 48	locs: 46	$\text{loc}(q_i) < \text{loc}(q_{i+1})$
98	49	51	
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	
		123	

© 2021, Jamie Callan

123

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	Do a left-to-right check
tf: 3	tf: 4	tf: 4	of locations
locs: 47	locs: 48	locs: 46	• $48 - 47 \leq n$ (match)
98	49	51	• $51 - 48 \leq n$ (match)
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	
		124	

© 2021, Jamie Callan

124

Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	Increment all loc iterators
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	
	
		125	

© 2021, Jamie Callan

125

Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	Advance loc iterators $q_{i>0}$ left-to-right such that $loc(q_i) < loc(q_{i+1})$
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	
	
		126	

© 2021, Jamie Callan

126

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	Do a left-to-right check of locations
tf: 3	tf: 4	tf: 4	
locs: 47	locs: 48	locs: 46	
98	49	51	133 - 98 > n (no match)
132	133	114	
doc: 92	134	137	
...	doc: 148	doc: 129	
	
		127	

© 2021, Jamie Callan

127

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	Increment q_0 loc iterator
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	
	
		128	

© 2021, Jamie Callan

128

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

		129

Query: #NEAR/3 (a b c)

Advance loc iterators $q_{i>0}$
left-to-right such that
 $\text{loc}(q_i) < \text{loc}(q_{i+1})$

- Not necessary here

© 2021, Jamie Callan

129

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

		130

Query: #NEAR/3 (a b c)

Do a left-to-right check
of locations

$133 - 132 \leq n$ (match)

$137 - 133 > n$ (no match)

© 2021, Jamie Callan

130

Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	Increment q_0 loc iterator
doc: 19	doc: 23	doc: 23	q_0 locs are exhausted.
tf: 1	tf: 1	tf: 1	No more matches are
locs: 7	locs: 99	loc: 99	possible in this document
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	
	
		131	

© 2021, Jamie Callan

131

Proximity Operators: A Simple Greedy NEAR/n Algorithm

a	b	c	Query: #NEAR/3 (a b c)
df: 47	df: 95	df: 14	Increment all doc iterators
doc: 19	doc: 23	doc: 23	...
tf: 1	tf: 1	tf: 1	Continue until the inverted
locs: 7	locs: 99	loc: 99	lists are exhausted
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	
	

© 2021, Jamie Callan

132

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

		133

Query: #NEAR/3 (a b c)

**Perhaps you expected q_0 's
loc iterator to be advanced
when this match failed**

**This is a flaw in the simple
greedy algorithm.**

**Better algorithms are
possible, but also more
complex**

© 2021, Jamie Callan

133

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 {0, 1} and {13, 15}
- Results for the NEAR operator: tf=2, and locations=1, 15

Query: #NEAR/2 (a b c)

Text: a a b b c c

- There are two matches {0, 2, 4} and {1, 3, 5}
- Results for the NEAR operator: tf=2, and locations=4, 5

134

© 2021, Jamie Callan

134

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

135

© 2021, Jamie Callan

135

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 {0, 1, 6} (the first location for each term)
 - {0, 1, 6} fails to match, so the q_0 location pointer advances
 - » a
 - The list for a is exhausted, so this text does not match
- A better algorithm would find a match at {0, 3, 6}
 - More accurate algorithms are much slower
 - The greedy algorithm is usually sufficient in practice
- Use the greedy algorithm for your homework

136

© 2021, Jamie Callan

136

Proximity Operators: NEAR/n FAQ

Query: #NEAR/2 (a a b)

Text: a b a a b

- Our algorithm is not explicitly designed for duplicate arguments
- But ... nothing prohibits duplicate arguments
- Probably it will work in most cases
 - E.g., it would match locations (0, 2, 4) above

These cases haven't been studied much by researchers

- But Google seems to support “apple apple pie”

137

© 2021, Jamie Callan

137

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”)

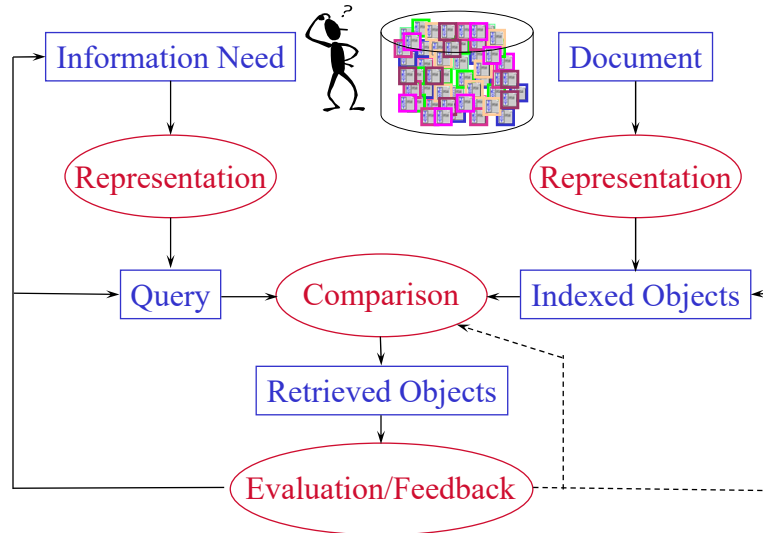
- Later lectures explore these topics in greater detail

138

© 2021, Jamie Callan

138

Overview of Information Retrieval Processes



139

© 2021, Jamie Cullin