

CS 50: Software Design and Implementation

Querier

Agenda



1. Querier

2. Fuzz testing

3. Activity

Crawler finds pages reachable from seedURL and stores URL, depth, HTML

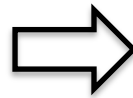
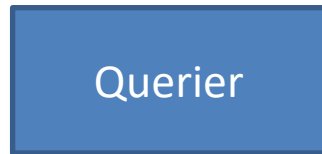
Goal:

- Keep track of to crawl pages
- Do not revisit pages

ADTs?

- Bag to track pages to see
- Hashtable for fast look up of pages seen

Query words



If your crawler didn't work well, find example output at:
</thayerfs/courses/22fall/cosc050/workspace/tse/tse-output>

Use these examples as a source for your indexer

Given:

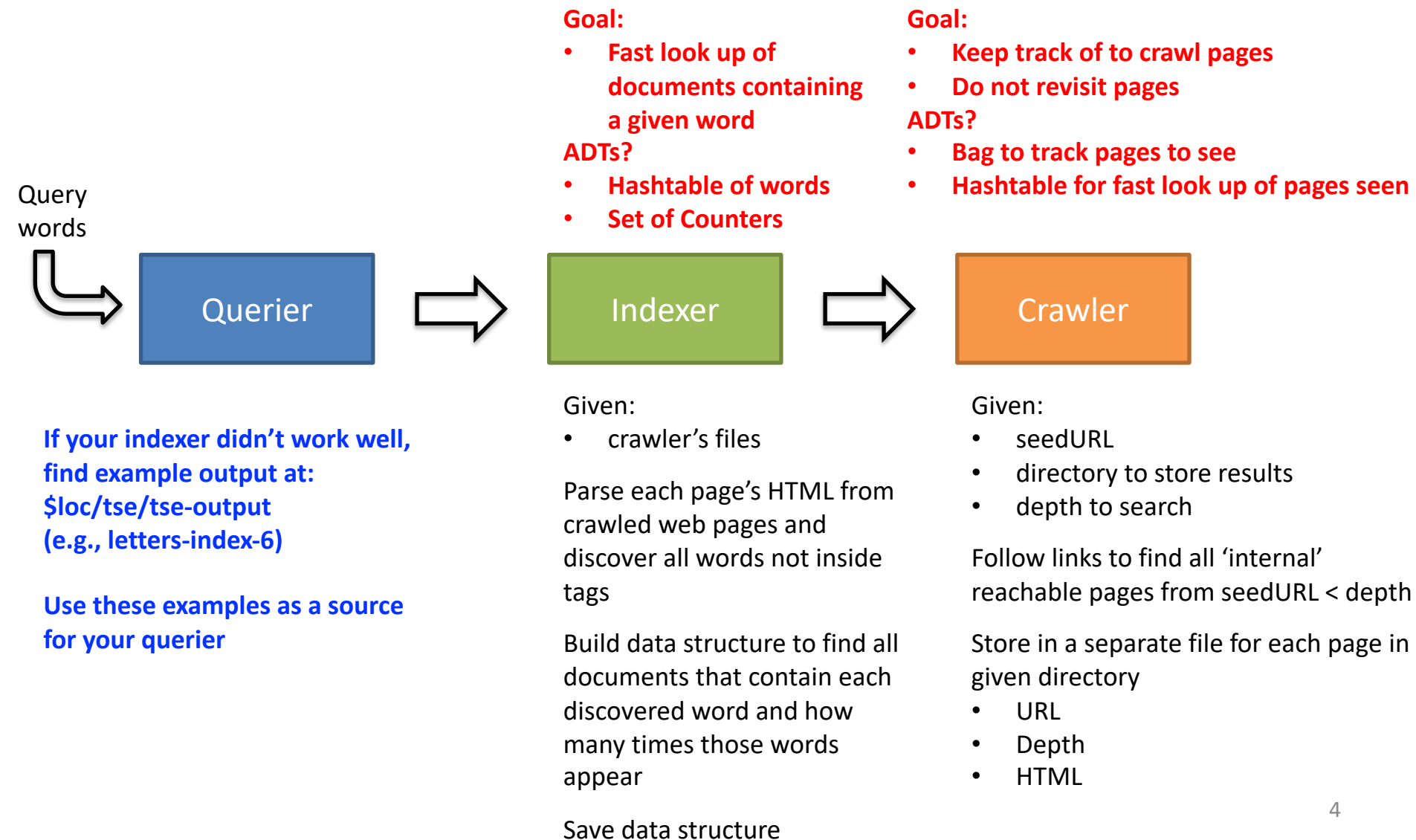
- seedURL
- directory to store results
- depth to search

Follow links to find all reachable pages from seedURL < depth

Store in a separate file for each page in given directory

- URL
- Depth
- HTML

Indexer uses crawler's results and builds data structure to find pages with words



Querier finds and ranks pages containing query words

Query words

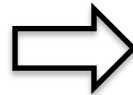
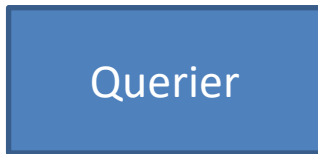


Goal:

- Fast ranked list of documents containing query words

ADTs?

- Indexer's hashtable of sets of counters



Goal:

- Fast look up of documents containing a given word

ADTs?

- Hashtable of words
- Set of Counters



Goal:

- Keep track of to crawl pages
- Do not revisit pages

ADTs?

- Bag to track pages to see
- Hashtable for fast look up of pages seen



Given:

- Indexer's data structure
- Query words

Find web pages containing query words

Rank pages based on how many times each word appears

Consider AND and OR logic

Given:

- crawler's files

Parse each page's HTML from crawled web pages and discover all words not inside tags

Build data structure to find all documents that contain each discovered word and how many times those words appear

Save data structure

Given:

- seedURL
- directory to store results
- depth to search

Follow links to find all 'internal' reachable pages from seedURL < depth

Store in a separate file for each page in given directory

- URL
- Depth
- HTML

Querier takes AND as well as OR queries

Query

computer science

- Implicit AND between computer and science
- Pages must have both words

computer and science

- Same as first query

computer or science

- Returns pages that have either computer OR science

baseball or basketball or ultimate frisbee

- Baseball OR basketball OR (ultimate AND frisbee)

DEMO

Enter search terms

```
$ ./querier $loc/tse/tse-output/letters-depth-6 $loc/tse/tse-output/letters-index-6
```

```
Query? first and search
```

Pass pages directory and index

```
Query: first and search
```

```
Matches 2 documents (ranked):
```

```
score 1 doc 3: http://cs50tse.cs.dartmouth.edu/tse/letters/B.html
```

```
score 1 doc 8: http://cs50tse.cs.dartmouth.edu/tse/letters/D.html
```

Print normalized query

```
Query? tiny search engine
```

Score for each page

Return results based on number of times word appears on each page

- AND: min of scores
- OR: add scores

```
Query: tiny search engine
```

```
No documents match.
```

If no matches, say so

```
Query? NOTE we LOWERcase the query
```

Page URL

```
Query: note we lowercase the query
```

```
No documents match.
```

Normalized changes query entered to lower case

```
Query? spaces do not matter
```

Ignore spaces

```
Query: spaces do not matter
```

```
No documents match.
```

DEMO

Query? **non-letter characters are disallowed**

Error: bad character '-' in query.

Query? **even digits as in cs50**

Error: bad character '5' in query.

Query? **and**

Query: and

Error: 'and' cannot be first

Query? **or**

Query: or

Error: 'or' cannot be first

Query? **what about and**

Query: what about and

Error: 'and' cannot be last

Query? **what about or**

Query: what about or

Error: 'or' cannot be last

Query? ^D

Only allow alphabetic characters

Queries cannot start or end with AND or OR

Queries also cannot have two operators (AND, OR) consecutively

Agenda

1. Querier

 2. Fuzz testing

3. Activity

Fuzz testing sends random input to the querier

```
202 static void
203 generateQuery(const wordlist_t* wordlist, const wordlist_t* dictionary)
204 {
205     // some parameters that affect query generation
206     const int maxWords = 6;           // generate 1..maxWords
207     const float orProbability = 0.3;  // P(OR between two words)
208     const float andProbability = 0.2; // P(AND between two words)
209     const float dictProbability = 0.2; // P(draw from dict instead of wordlist)
210
211     int qwords = rand() % maxWords + 1; // number of words in query
212     for (int qw = 0; qw < qwords; qw++) {
213         // draw a word either dictionary or wordlist
214         if ((rand() % 100) < (dictProbability * 100)) {
215             printf("%s ", dictionary->words[rand() % dictionary->nWords]);
216         } else {
217             printf("%s ", wordlist->words[rand() % wordlist->nWords]);
218         }
219
220         // last word?
221         if (qw < qwords-1) {
222             // which operator to print?
223             int op = rand() % 100;
224             if (op < (andProbability * 100)) {
225                 printf("AND ");
226             }
227             else if (op < (andProbability * 100 + orProbability * 100)) {
228                 printf("OR ");
229             }
230         }
231     }
232     printf("\n");
233 }
```

Words from index

Words from dictionary

Generate query with qwords in it

Pick random word from either words from index or dictionary

Put an AND or OR operator randomly, but not first or last!

Fuzz testing sends random input to the querier

```
202 static void
203 generateQuery(const wordlist_t* wordlist, const wordlist_t* dictionary)
204 {
205     // some parameters that affect query generation
206     const int maxWords = 6;           // generate 1..maxWords
207     const float orProbability = 0.3;  // P(OR between two words)
208     const float andProbability = 0.2; // P(AND between two words)
209     const float dictProbability = 0.2; // P(draw from dict instead of wordlist)
210
211     int qwords = rand() % maxWords + 1; // number of words in query
212     for (int qw = 0; qw < qwords; qw++) {
213         // draw a word either dictionary or wordlist
214         if ((rand() % 100) < (dictProbability * 100)) {
215             printf("%s ", dictionary->words[rand() % dictionary->wordCount]);
216         } else {
217             printf("%s ", wordlist->words[rand() % wordlist->wordCount]);
218         }
219
220         // last word?
221         if (qw < qwords-1) {
222             // which operator to print?
223             int op = rand() % 100;
224             if (op < (andProbability * 100)) {
225                 printf("AND ");
226             }
227             else if (op < (andProbability * 100 + orProbability * 100)) {
228                 printf("OR ");
229             }
230         }
231     }
232     printf("\n");
233 }
```

Words from index
Number of queries to make
Random seed

```
$. /fuzzquery $loc/tse/tse-output/letters-index-6 10 0
./fuzzquery: generating 10 queries from 22 words
fourier AND traversal
this OR the the OR tse computational
biology playground OR computational
answers breadth search OR computational OR Mississippians
OR fast
algorithm OR coding eniac the AND home OR breadth
traversal computational playground coding OR the
fast
search the OR fast
home
transform OR huffman OR depth AND graph AND transform
```

Pipe fuzzer's queries to querier to test

```
$ ./fuzzquery $loc/tse/tse-output/letters-index-6 10 0 | ./querier $loc/tse/tse-output/letters-depth-6 $loc/tse/tse-output/letters-index-6
```

```
./fuzzquery: generating 10 queries from 22 words  
Query: fourier and traversal  
No documents match.
```

Crawled pages **Indexed pages**



```
-----  
Query: this or the the or tse computational  
Matches 1 documents (ranked):  
score 2 doc 1: http://cs50tse.cs.dartmouth.edu/tse/letters/index.html
```

```
-----  
Query: biology playground or computational  
Matches 1 documents (ranked):  
score 1 doc 9: http://cs50tse.cs.dartmouth.edu/tse/letters/C.html
```

```
-----  
Query: answers breadth search or computational or mississippians or fast  
Matches 2 documents (ranked):  
score 1 doc 9: http://cs50tse.cs.dartmouth.edu/tse/letters/C.html  
score 1 doc 7: http://cs50tse.cs.dartmouth.edu/tse/letters/F.html
```

```
-----  
Query: algorithm or coding eniac the and home or breadth  
Matches 2 documents (ranked):  
score 1 doc 2: http://cs50tse.cs.dartmouth.edu/tse/letters/A.html  
score 1 doc 3: http://cs50tse.cs.dartmouth.edu/tse/letters/B.html
```

```
-----  
Query: traversal computational playground coding or the  
Matches 1 documents (ranked):  
score 1 doc 1: http://cs50tse.cs.dartmouth.edu/tse/letters/index.html
```

```
-----  
Query: fast  
Matches 1 documents (ranked):  
score 1 doc 7: http://cs50tse.cs.dartmouth.edu/tse/letters/F.html
```


Issues?

- **Should change random seed (otherwise queries always the same)**
- **Does not check if results are correct**
- **Does show if program crashes!**

Agenda

1. Querier

2. Fuzz testing

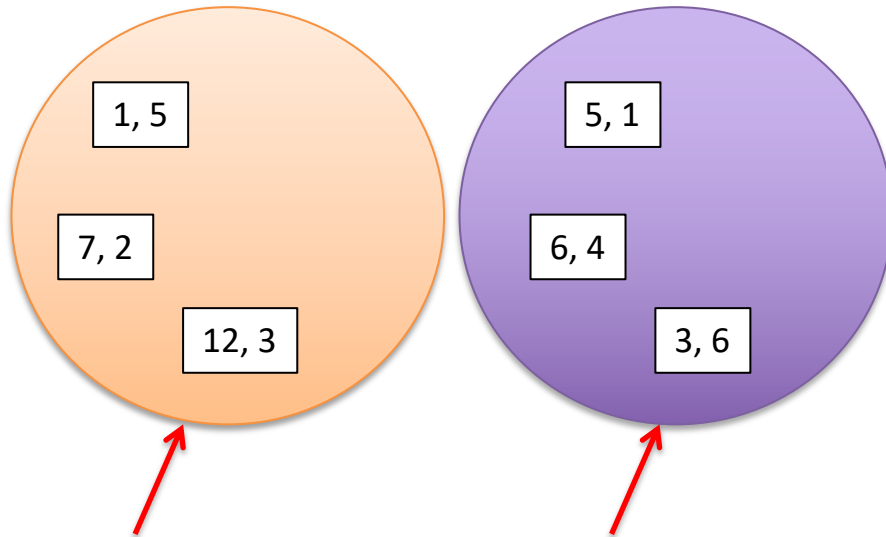
 3. Activity

Users enter query words, and the querier must implement AND and OR operations

Note: currently no pages contain both words

Query word1

Query word2



DocID, count of pages containing dartmouth

DocID, count of pages containing algorithm

Example:
Doc1 contains word1 5 times

Example:
Doc 6 contains word2 4 times

Example

Word 1 = dartmouth
word 2 = algorithm

Dartmouth appears on sites 1, 7, and 12

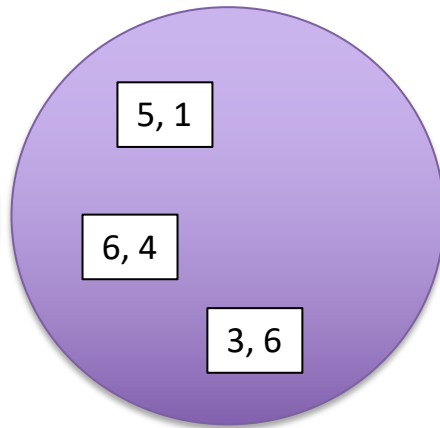
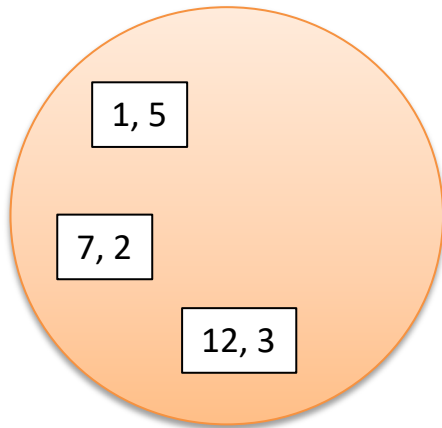
Algorithm appears on sites 3, 5, and 6

OR is the UNION of two sets

Note: currently no pages contain both words

Query word1

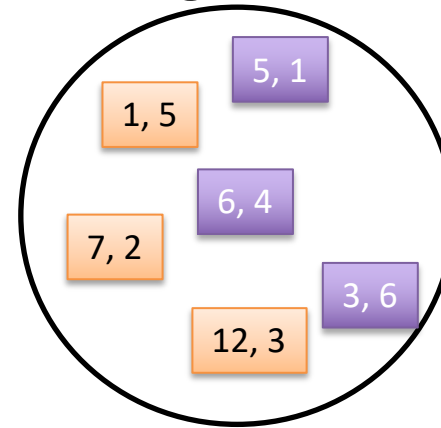
Query word2



dartmouth

OR

algorithm



OR returns sites that mention either word (UNION)

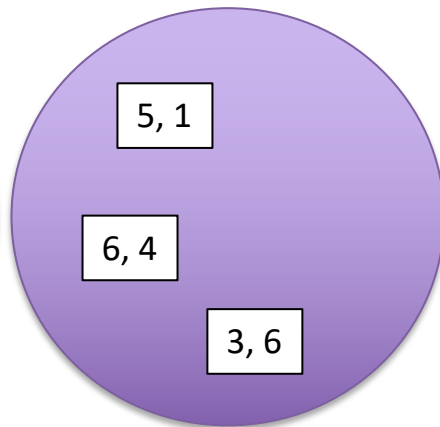
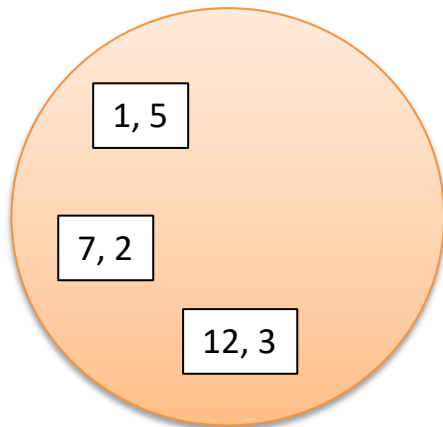
Return sites sorted by count
(3, 1, 6, 12, 7, 5)

AND is the INTERSECTION of two sets

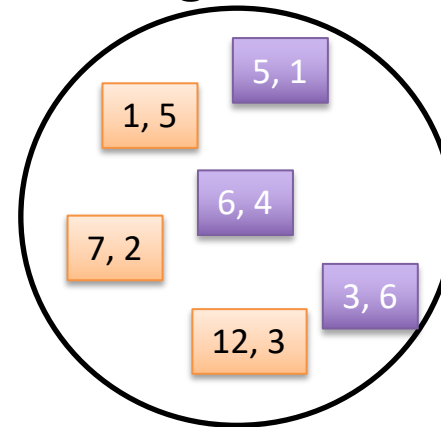
Note: currently no pages contain both words

Query word1

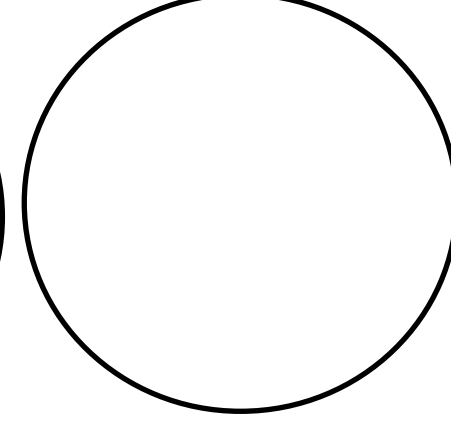
Query word2



dartmouth
OR
algorithm



dartmouth
AND
algorithm



OR returns sites that mention either word (UNION)

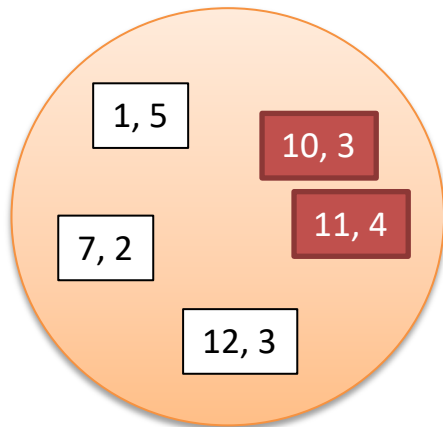
AND returns sites that mention both words (INTERSECTION)

Return sites sorted by count
(3, 1, 6, 12, 7, 5)

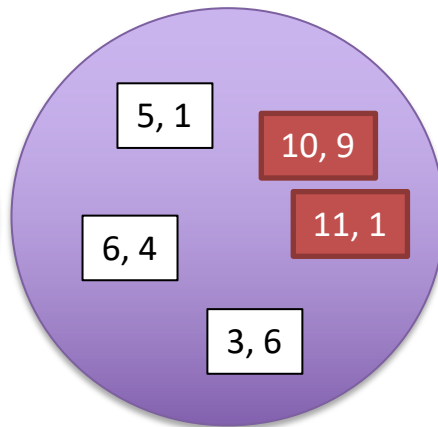
Currently none

Sometimes the same site contains multiple query words

Query word1



Query word2



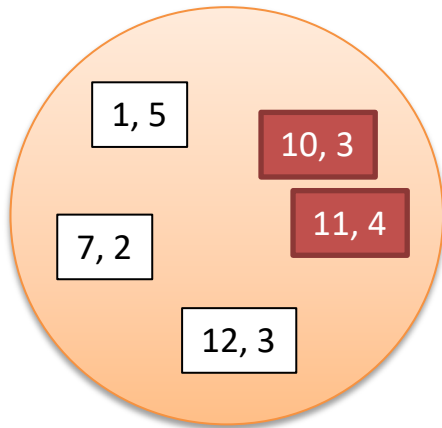
Page 10 contains Dartmouth 3 times and algorithm 9 times

Page 11 contains Dartmouth 4 times and algorithm 1 time

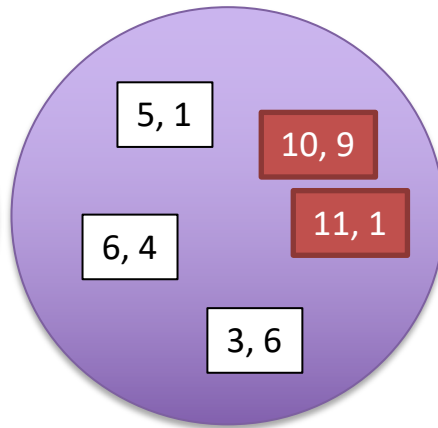
Now suppose pages 10 and 11 each contain query word 1 and 2

OR adds the counts from each site

Query word1



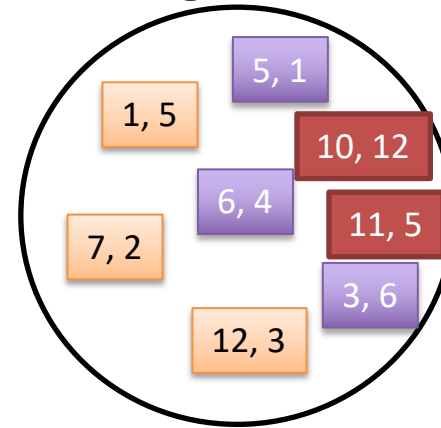
Query word2



dartmouth

OR

algorithm



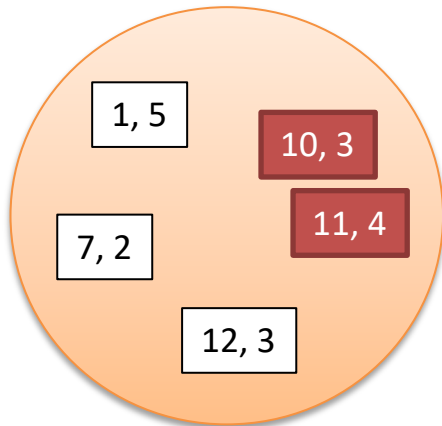
Now suppose pages 10 and 11
each contain query word 1 and 2

OR adds counts from
each site for overlaps

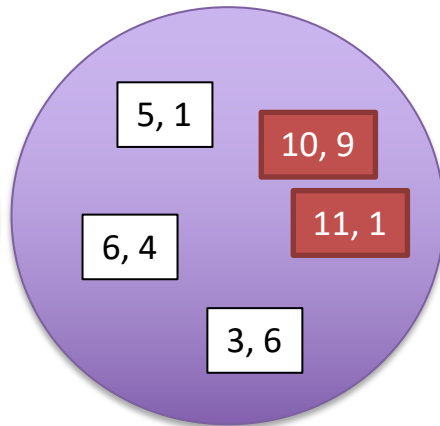
Still return sites sorted
by count
(10, 3, 11, 1, 6, 12, 7, 5)

AND takes the minimum count between both sites

Query word1

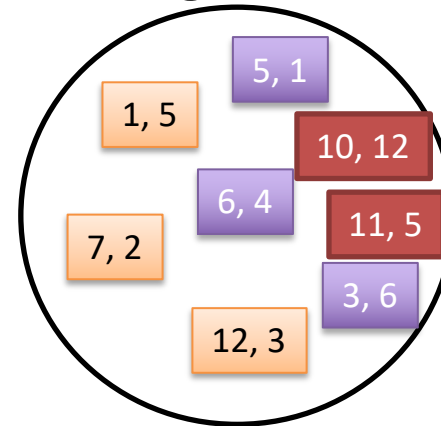


Query word2



Now suppose pages 10 and 11 each contain query word 1 and 2

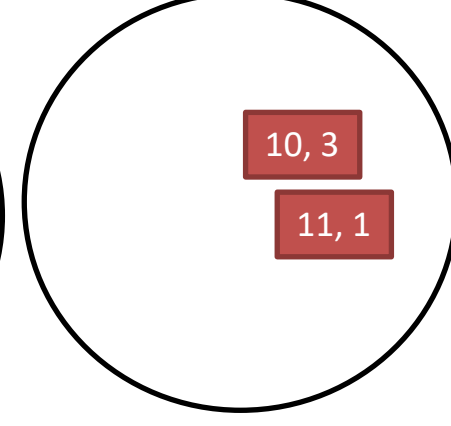
dartmouth
OR
algorithm



OR adds counts from each site for overlaps

Still return sites sorted by count
(10, 3, 11, 1, 6, 12, 7, 5)

dartmouth
AND
algorithm

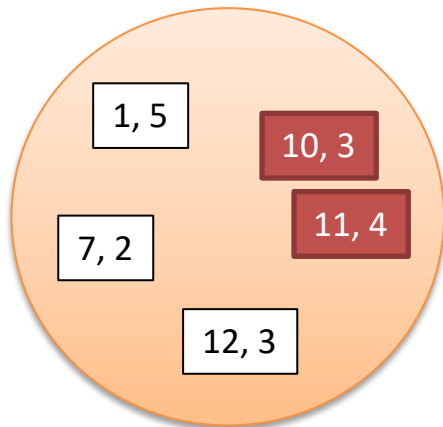


AND takes the minimum of each site

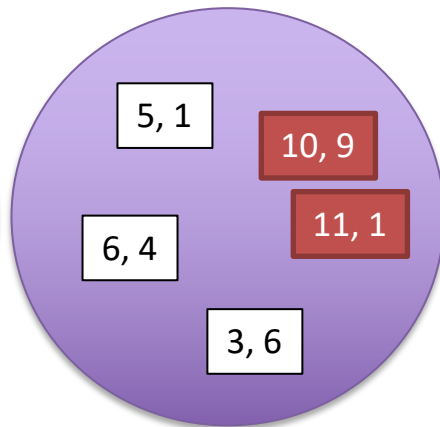
Return sites sorted by count
(10, 11)

AND takes the minimum count between both sites

Query word1



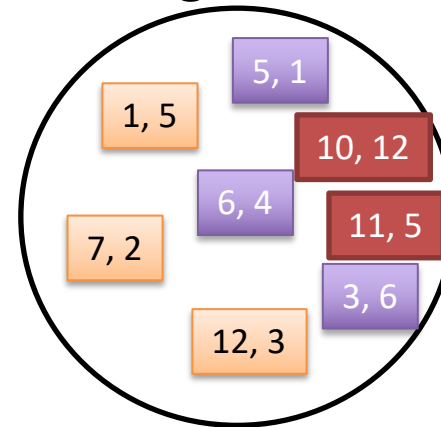
Query word2



Now suppose pages 10 and 11 each contain query word 1 and 2

If query words appear on more than two sites, add counts for OR, take min for AND

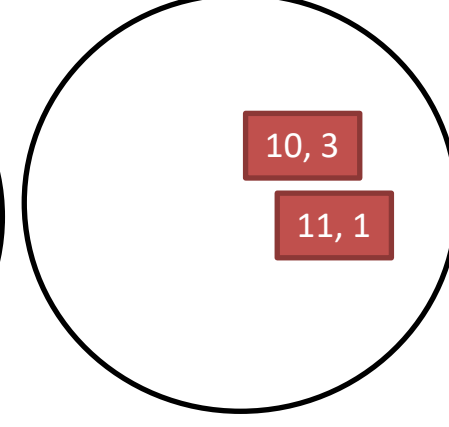
dartmouth
OR
algorithm



OR adds counts from each site for overlaps

Still return sites sorted by count
(10, 3, 11, 1, 6, 12, 7, 5)

dartmouth
AND
algorithm



AND takes the minimum of each site

Return sites sorted by count
(10, 11)

set_iterate2.c demonstrates UNION

set_iterate2.c

```
20 int main()
21 {
22     set_t *setA, *setB, *result; // three sets
```

Create three sets

```
23
24     setA = mem_assert(set_new(), "setA");
25     setB = mem_assert(set_new(), "setB");
26     result = mem_assert(set_new(), "result");
27
```

```
28     printf("Empty result set: ");
29     set_print(result, stdout, itemprint);
30     putchar('\n');
```

```
31
32     printf("Building set A: ");
33     set_insert(setA, "Brown", intsave(10));
34     set_insert(setA, "Dartmouth", intsave(20));
35     set_insert(setA, "Yale", intsave(15));
36     set_insert(setA, "Harvard", intsave(8));
37     set_insert(setA, "Princeton", intsave(5));
38     set_insert(setA, "Columbia", intsave(1));
39     set_print(setA, stdout, itemprint);
40     putchar('\n');
```

Load setA and setB with college string as key and a score as item

```
41
42     printf("Building set B: ");
43     set_insert(setB, "Penn", intsave(7));
44     set_insert(setB, "Dartmouth", intsave(11));
45     set_insert(setB, "Cornell", intsave(9));
46     set_insert(setB, "Stanford", intsave(6));
47     set_insert(setB, "Princeton", intsave(3));
48     set_insert(setB, "Duke", intsave(12));
49     set_print(setB, stdout, itemprint);
```

```
114 static int *
115 intsave(int item)
116 {
117     int *saved =
118         mem_assert(malloc(sizeof(int)), "intsave");
119     *saved = item;
120 }
```

Why do we call intsave?
set_insert takes a pointer to an item

Merge setA into result, then merge setB into result

```
52 printf("\nMerge of setA into result: \n");
53 set_merge(result, setA);
54 set_print(result, stdout, itemprint); ←
55 putchar('\n');
56
57 printf("\nMerge of setB into result: \n");
58 set_merge(result, setB);
59 set_print(result, stdout, itemprint); ←
60 putchar('\n');
61
62 printf("\nDelete the sets...\n");
63 set_delete(setA, itemdelete);
64 set_delete(setB, itemdelete);
65 set_delete(result, itemdelete); ←
66 }
```

Result is initially empty set_iterate2.c

**Merge (UNION) setA into result
Print result**

**Merge setB into result
Print result**

Delete all three sets when done

set_merge adds items from two sets if the key is in both sets, else insert key, item

```
68 /* Merge the second set into the first set;
69 * the second set is unchanged.
70 */
```

```
71 static void
72 set_merge(set_t *setA, set_t *setB)
73 {
74     set_iterate(setB, setA, set_merge_helper);
75 }
```

```
76
77 /* Consider one item for insertion into the other set.
78 * If the other set does not contain the item, insert it;
79 * otherwise, update the other set's item with sum of item values.
80 */
```

```
81 static void
82 set_merge_helper(void *arg, const char *key, void *item)
83 {
```

```
84     set_t *setA = arg;
85     int *itemB = item;
86
87     // find the same key in setA
88     int *itemA = set_find(setA, key);
89     if (itemA == NULL) {
90         // not found: insert it
91         set_insert(setA, key, intsave(*itemB));
92         printf("\t%s added\n", key);
93     } else {
94         // add to the existing value
95         *itemA += *itemB;
96         printf("\t%s exists\n", key);
97     }
98 }
```

set_iterate2.c

set_merge does a UNION of two sets

- Saves result in first set
- set_merge_helper called on each item in second set

**setA passed as arg
Cast as set_t
Cast item as integer pointer to
setB's item**

**If setA does not have key from
setB, add key and integer to setA**

**Otherwise, add itemB to itemA
to update setA**

Activity is to implement INTERSECT given starter code in counters_intersect.c

```
20 int main()
21 {
22     // create two counters for demo
23     counters_t *c1 = mem_assert(counters_new(), "counters_new() failed");
24     counters_t *c2 = mem_assert(counters_new(), "counters_new() failed");
25
26     // init counters 1
27     counters_set(c1, 3, 6);
28     counters_set(c1, 4, 7);
29     counters_set(c1, 5, 1);
30     counters_set(c1, 7, 4);
31
32     // init counters 2
33     counters_set(c2, 1, 3);
34     counters_set(c2, 3, 2);
35     counters_set(c2, 5, 4);
36     counters_set(c2, 6, 6);
37     counters_set(c2, 7, 3);
38
39     // take the intersection, store the results in c1
40     counters_intersect(c1, c2);
41     counters_print(c1, stdout);
42     printf("\n");
43
44     // clean up
45     counters_delete(c1);
46     counters_delete(c2);
47
48     return 0;
49 }
50
51 // TODO: fill in this function
52 void counters_intersect(counters_t* ct1, counters_t* ct2)
53 {
54
55 }
```

Complete counters_intersect to keep minimum of counts where keys match

