

# Chapter IR:II

## II. Indexing

- ❑ Indexing Basics
- ❑ Inverted Index
- ❑ Query Processing I
- ❑ Query Processing II
- ❑ Index Construction
- ❑ Index Compression
- ❑ Size Estimation

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$\dots$
$t_1$						
$t_2$						
$t_3$						
$t_4$						
$t_5$						
$\vdots$						$\dots$

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...
$t_1$						
$t_2$						
$t_3$						
$t_4$						
$t_5$						
⋮						...

### □ Documents $D$

$d_1$  Antony and Cleopatra

$d_2$  Julius Caesar

$d_3$  The Tempest

$d_4$  Hamlet

$d_5$  Othello

### □ Index terms $T$

$t_1$  Antony

$t_2$  Brutus

$t_3$  Caesar

$t_4$  Calpurnia

$t_5$  Cleopatra

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...
$t_1$	1					
$t_2$	1					
$t_3$	1					
$t_4$	0					
$t_5$	1					
⋮						...

### □ Documents $D$

$d_1$  Antony and Cleopatra

$d_2$  Julius Caesar

$d_3$  The Tempest

$d_4$  Hamlet

$d_5$  Othello

### □ Index terms $T$

$t_1$  Antony

$t_2$  Brutus

$t_3$  Caesar

$t_4$  Calpurnia

$t_5$  Cleopatra

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...
$t_1$	1	1	0	0	0	
$t_2$	1	1	0	1	0	
$t_3$	1	1	0	1	1	
$t_4$	0	1	0	0	0	
$t_5$	1	0	0	0	0	
⋮						...

### □ Documents $D$

$d_1$  Antony and Cleopatra

$d_2$  Julius Caesar

$d_3$  The Tempest

$d_4$  Hamlet

$d_5$  Othello

### □ Index terms $T$

$t_1$  Antony

$t_2$  Brutus

$t_3$  Caesar

$t_4$  Calpurnia

$t_5$  Cleopatra

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...
$t_1$	382	128	0	0	0	
$t_2$	4	379	0	1	0	
$t_3$	289	272	0	2	1	
$t_4$	0	16	0	0	0	
$t_5$	271	0	0	0	0	
⋮						...

### □ Documents $D$

$d_1$  [Antony and Cleopatra](#)

$d_2$  [Julius Caesar](#)

$d_3$  [The Tempest](#)

$d_4$  [Hamlet](#)

$d_5$  [Othello](#)

### □ Index terms $T$

$t_1$  Antony

$t_2$  Brutus

$t_3$  Caesar

$t_4$  Calpurnia

$t_5$  Cleopatra

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...
$t_1$	$w_{1,1}$	$w_{1,2}$	$w_{1,3}$	$w_{1,4}$	$w_{1,5}$	
$t_2$	$w_{2,1}$	$w_{2,2}$	$w_{2,3}$	$w_{2,4}$	$w_{2,5}$	
$t_3$	$w_{3,1}$	$w_{3,2}$	$w_{3,3}$	$w_{3,4}$	$w_{3,5}$	
$t_4$	$w_{4,1}$	$w_{4,2}$	$w_{4,3}$	$w_{4,4}$	$w_{4,5}$	
$t_5$	$w_{5,1}$	$w_{5,2}$	$w_{5,3}$	$w_{5,4}$	$w_{5,5}$	
⋮						...

### □ Documents $D$

$d_1$  Antony and Cleopatra

$d_2$  Julius Caesar

$d_3$  The Tempest

$d_4$  Hamlet

$d_5$  Othello

### □ Index terms $T$

$t_1$  Antony

$t_2$  Brutus

$t_3$  Caesar

$t_4$  Calpurnia

$t_5$  Cleopatra

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...
$t_1$	$w_{1,1}$	$w_{1,2}$	$w_{1,3}$	$w_{1,4}$	$w_{1,5}$	
$t_2$	$w_{2,1}$	$w_{2,2}$	$w_{2,3}$	$w_{2,4}$	$w_{2,5}$	
$t_3$	$w_{3,1}$	$w_{3,2}$	$w_{3,3}$	$w_{3,4}$	$w_{3,5}$	
$t_4$	$w_{4,1}$	$w_{4,2}$	$w_{4,3}$	$w_{4,4}$	$w_{4,5}$	
$t_5$	$w_{5,1}$	$w_{5,2}$	$w_{5,3}$	$w_{5,4}$	$w_{5,5}$	
⋮						...

### □ Documents $D$

- $d_1$  Antony and Cleopatra
- $d_2$  Julius Caesar
- $d_3$  The Tempest
- $d_4$  Hamlet
- $d_5$  Othello

### □ Index terms $T$

- $t_1$  Antony
- $t_2$  Brutus
- $t_3$  Caesar
- $t_4$  Calpurnia
- $t_5$  Cleopatra

### □ Term Weights

- Boolean
- Term frequency
- ...



# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$\dots$
$t_1$	$w_{1,1}$	$w_{1,2}$	$w_{1,3}$	$w_{1,4}$	$w_{1,5}$	
$t_2$	$w_{2,1}$	$w_{2,2}$	$w_{2,3}$	$w_{2,4}$	$w_{2,5}$	
$t_3$	$w_{3,1}$	$w_{3,2}$	$w_{3,3}$	$w_{3,4}$	$w_{3,5}$	
$t_4$	$w_{4,1}$	$w_{4,2}$	$w_{4,3}$	$w_{4,4}$	$w_{4,5}$	
$t_5$	$w_{5,1}$	$w_{5,2}$	$w_{5,3}$	$w_{5,4}$	$w_{5,5}$	
$\vdots$						$\dots$

### Observations:

- ❑ Most retrieval models induce a term-document matrix by computing term weights  $w_{i,j}$  for each pair of term  $t_i \in T$  and document  $d_j \in D$ .
- ❑ Query-independent computations that depend only on  $D$  are done offline.
- ❑ Online, for a query  $q$ , the required term weights are looked up to score documents.

# Inverted Index

## Term-Document Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$\dots$
$t_1$	$w_{1,1}$	$w_{1,2}$				
$t_2$	$w_{2,1}$	$w_{2,2}$		$w_{2,4}$		
$t_3$	$w_{3,1}$	$w_{3,2}$		$w_{3,4}$	$w_{3,5}$	
$t_4$		$w_{4,2}$				
$t_5$	$w_{5,1}$					
$\vdots$						$\dots$

### Observations:

- ❑ The size of the term-document matrix is  $|T| \cdot |D|$ .
- ❑ The term-document matrix is sparse: the vast majority of term weights are 0.
- ❑ Therefore, most of the storage space required for the full matrix is wasted.
- ❑ Using a sparse-matrix representation yields significant space savings.
- ➔ An inverted index efficiently encodes a sparse term-document matrix.

# Inverted Index

## Data Structure

---

$T$	→	<b>Postings</b> (Posting Lists, Postlists)			
$t_1$	→	$d_1, w_{1,1}$	$d_2, w_{1,2}$		
$t_2$	→	$d_1, w_{2,1}$	$d_2, w_{2,2}$	$d_4, w_{2,4}$	
$t_3$	→	$d_1, w_{3,1}$	$d_2, w_{3,2}$	$d_4, w_{3,4}$	$d_5, w_{3,5}$
$t_4$	→	$d_2, w_{4,2}$			
$t_5$	→	$d_1, w_{5,1}$			
$\vdots$					

---

An index is implemented as a multimap (i.e., a hash table with multiple values).

Components of an externalized implementation:

- ❑ Term vocabulary file

Lookup table which maps terms  $t_i \in T$  to the start of their posting list in the postings file.

- ❑ Postings file(s)

File(s) that store posting lists on disk.

- ❑ Index entries  $d_i, [\dots]$ , so-called postings

# Inverted Index

## Data Structure

$T$	→	<b>Postings</b> (Posting Lists, Postlists)			
$t_1$	→	$d_1, w_{1,1}$	$d_2, w_{1,2}$		
$t_2$	→	$d_1, w_{2,1}$	$d_2, w_{2,2}$	$d_4, w_{2,4}$	
$t_3$	→	$d_1, w_{3,1}$	$d_2, w_{3,2}$	$d_4, w_{3,4}$	$d_5, w_{3,5}$
$t_4$	→	$d_2, w_{4,2}$			
$t_5$	→	$d_1, w_{5,1}$			
$\vdots$					

An index is implemented as a multimap (i.e., a hash table with multiple values).

Design choices:

- ❑ Information stored in a posting  $d_i, [\dots]$ .
- ❑ Ordering of each term's posting list.
- ❑ Encoding and compression techniques for further space savings.
- ❑ Physical implementation details, such as external memory and distribution.

# Inverted Index

## Posting

Given term  $t$  and document  $d$ , their posting may include the following:

`<document> [ <weights> ] [ <positions> ] ...`

`<document>`:

- ❑ Reference to the document  $d$  in which term  $t$  occurs (or to which it applies).

`<weights>`:

- ❑ Term weight  $w$  for term  $t$  in document  $d$ .
- ❑ Often, only basic term weights are stored (e.g., term frequency  $tf(t, d)$ ).  
Storing model-specific weights saves runtime at the expense of flexibility.

`<positions>`:

- ❑ Term positions within the document, e.g., term, sentence, page, chapter, etc.
- ❑ Field information, e.g., title, author, introduction, etc.

# Inverted Index

## Posting

Two special-purpose entries are distinguished:

... [`<list length>`]

... [`<skip pointer>`]

`<list length>`:

- ❑ Added to the first entry of the posting list of a term  $t$ .
- ❑ Stores the length of the posting list.
- ❑ What does the length of a posting list indicate?

`<skip pointer>`:

- ❑ Used to implement a [skip list](#) in a term's posting list, when ordered by ID.
- ❑ Allows for random access to postings in  $O(\log df(t, D))$ .
- ❑ An effective amount of skip entries has been found to be  $\sqrt{df(t, D)}$ .  
First entry of a posting list, and then at random (or regular) intervals.

# Inverted Index

## Posting

Two special-purpose entries are distinguished:

... [`<list length>`]

... [`<skip pointer>`]

`<list length>`:

- ❑ Added to the first entry of the posting list of a term  $t$ .
- ❑ Stores the length of the posting list.
- ❑ Equals the number of documents containing  $t$  (document frequency  $df(t, D)$ ).

`<skip pointer>`:

- ❑ Used to implement a [skip list](#) in a term's posting list, when ordered by ID.
- ❑ Allows for random access to postings in  $O(\log df(t, D))$ .
- ❑ An effective amount of skip entries has been found to be  $\sqrt{df(t, D)}$ .  
First entry of a posting list, and then at random (or regular) intervals.

# Inverted Index

## Posting List, Postlist

Example for two posting lists, where for term  $t_i$  postings  $[k, tf(t_i, d_k)]$  are stored:

$T$	Postings															
$\vdots$																
$t_i$	2, 4			4, 9	8, 2	16, 1		19, 7	23, 5	28, 6		41, 8	50, 6	77, 8		...
$t_j$	1, 1			2, 3	3, 5	5, 2		8, 17	41, 6	51, 5		60, 5	71, 3	77, 2		...
$\vdots$																

Ordering:

- ❑ by document identifier. Problem: “good” documents randomly distributed.
- ❑ by document quality. Problem: index updates more complicated.
- ❑ by term weight. Problem: no canonical order across rows; skip lists useless.

Compression:

- ❑ The size of an index is in  $O(|D|)$ , where  $|D|$  denotes the disk size of  $D$ .
- ❑ Posting lists can be effectively compressed with tailored techniques.



## Remarks:

- ❑ The term “inverted index” is redundant: “index” already denotes the structure in which terms are assigned to the (parts of) documents in which they occur. Better suited, but less frequently used, is “inverted file”, which expresses that a (document) file is “inverted” to form an index. So instead of assigning terms to documents, an index assigns documents to terms.
- ❑ A trade-off must be made between the amount of information stored in a posting and the time required to process a post list. The more information stored in a posting, the more has to be loaded into memory and decoded as the posting list is traversed.
- ❑ A skip entry can contain more than one pointer, so skip steps of different lengths are possible.
- ❑ Depending on the search domain, it may be beneficial to create more than one index with different properties.

# Chapter IR:II

## II. Indexing

- ❑ Indexing Basics
- ❑ Inverted Index
- ❑ Query Processing I
- ❑ Query Processing II
- ❑ Index Construction
- ❑ Index Compression
- ❑ Size Estimation

# Query Processing I

## Retrieval Types

Query processing can be based on two basic approaches:

- ❑ **Set retrieval**

A query induces a subset of the indexed documents which is considered relevant.  
Important applications: e-discovery, patent search, systematic reviews.

- ❑ **Ranked retrieval**

A query induces a ranking among all indexed documents in descending order of relevance.

Ranked retrieval is the norm in virtually all modern search engines.

# Query Processing I

## Query Semantics for Set Retrieval

Keyword queries have Boolean semantics that is either implicitly specified by user behavior and expectations or explicitly specified.

We distinguish four types:

- ❑ Single-term queries
- ❑ Disjunctive multi-term queries  
Only Boolean OR connectives. Example:  $\text{Antony} \vee \text{Brutus} \vee \text{Calpurnia}$ .
- ❑ Conjunctive multi-term queries  
Only Boolean AND connectives. Example:  $\text{Antony} \wedge \text{Brutus} \wedge \text{Calpurnia}$ .
  - + Constraint: Proximity  
Example:  $\text{Antony} /5 \text{ Caesar}$
  - + Constraint: Phrase  
Example: “Antony and Caesar”
- ❑ “Complex” Boolean multi-term queries  
Remainder of Boolean formulas. Example:  $(\text{Antony} \vee \text{Caesar}) \wedge \neg \text{Calpurnia}$ .  
Normalized to disjunctive or conjunctive normal form.

## Remarks:

### ❑ Which index configuration applies to which type of query?

#### Query types:

- Single-term queries
- Disjunctive multi-term queries
- Conjunctive multi-term queries
  - Boolean AND queries
  - Proximity queries
  - Phrase queries

#### Index configurations:

- Postlists ordered by document ID
- Postlists ordered by document quality
- Postlists ordered by term weight
- Positional indexing  
Postings also store term positions.

## Remarks:

- ❑ Which index configuration applies to which type of query?

### Query types:

- Single-term queries
- Disjunctive multi-term queries
- Conjunctive multi-term queries
  - Boolean AND queries
  - Proximity queries
  - Phrase queries

### Index configurations:

- Postlists ordered by document ID
- Postlists ordered by document quality
- Postlists ordered by term weight
- Positional indexing  
Postings also store term positions.

- ❑ Single-term queries are directly answered with a term weight ordering.
- ❑ Disjunctive multi-term queries can be processed with any postlist ordering.
- ❑ Conjunctive multi-term queries benefit from a canonical postlist order.
- ❑ Proximity and phrase queries require positional indexing.

# Query Processing I

## Conjunctive Multi-Term Queries

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$  and a query  $q = t_1 \wedge \dots \wedge t_n$ , compute the collection  $R \subseteq D$  of documents relevant to  $q$ .

$T$	Postings
$\vdots$	
$t_i$	<div>2, <span style="color:blue">4</span> <span style="background-color:yellow; border:1px solid black; display:inline-block; width:10px; height:10px;"></span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span></div> <div>4, <span style="color:blue">9</span></div> <div>8, <span style="color:blue">2</span></div> <div>16, <span style="color:blue">1</span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span></div> <div>19, <span style="color:blue">7</span></div> <div>23, <span style="color:blue">5</span></div> <div>28, <span style="color:blue">6</span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span></div> <div>41, <span style="color:blue">8</span></div> <div>50, <span style="color:blue">6</span></div> <div>77, <span style="color:blue">8</span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span> ...</div>
$t_j$	<div>1, <span style="color:blue">1</span> <span style="background-color:yellow; border:1px solid black; display:inline-block; width:10px; height:10px;"></span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span></div> <div>2, <span style="color:blue">3</span></div> <div>3, <span style="color:blue">5</span></div> <div>5, <span style="color:blue">2</span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span></div> <div>8, <span style="color:blue">17</span></div> <div>41, <span style="color:blue">6</span></div> <div>51, <span style="color:blue">5</span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span></div> <div>60, <span style="color:blue">5</span></div> <div>71, <span style="color:blue">3</span></div> <div>77, <span style="color:blue">2</span> <span style="background-color:blue; border:1px solid black; display:inline-block; width:10px; height:10px;"></span> ...</div>
$\vdots$	

What is the underlying problem to which processing query  $q$  can be reduced?

# Query Processing I

## Conjunctive Multi-Term Queries

Given an index with postings  $[k, \textcolor{violet}{tf}(t, d_k)]$  and a query  $q = t_1 \wedge \dots \wedge t_n$ , compute the collection  $R \subseteq D$  of documents relevant to  $q$ .

$T$	Postings
$\vdots$	
$t_i$	<div>2, <span style="background-color: #ff8c00;">4</span> <span style="background-color: #38a83d;"> </span></div> <div>4, <span style="color: #800080;">9</span> <span style="background-color: #38a83d;"> </span></div> <div>8, <span style="color: #800080;">2</span> <span style="background-color: #38a83d;"> </span></div> <div>16, <span style="color: #800080;">1</span> <span style="background-color: #38a83d;"> </span></div> <div>19, <span style="color: #800080;">7</span> <span style="background-color: #38a83d;"> </span></div> <div>23, <span style="color: #800080;">5</span> <span style="background-color: #38a83d;"> </span></div> <div>28, <span style="color: #800080;">6</span> <span style="background-color: #38a83d;"> </span></div> <div>41, <span style="color: #800080;">8</span> <span style="background-color: #38a83d;"> </span></div> <div>50, <span style="color: #800080;">6</span> <span style="background-color: #38a83d;"> </span></div> <div>77, <span style="color: #800080;">8</span> <span style="background-color: #38a83d;"> </span></div> <div>...</div>
$t_j$	<div>1, <span style="color: #800080;">1</span> <span style="background-color: #ff8c00;"> </span> <span style="background-color: #38a83d;"> </span></div> <div>2, <span style="color: #800080;">3</span> <span style="background-color: #38a83d;"> </span></div> <div>3, <span style="color: #800080;">5</span> <span style="background-color: #38a83d;"> </span></div> <div>5, <span style="color: #800080;">2</span> <span style="background-color: #38a83d;"> </span></div> <div>8, <span style="color: #800080;">17</span> <span style="background-color: #38a83d;"> </span></div> <div>41, <span style="color: #800080;">6</span> <span style="background-color: #38a83d;"> </span></div> <div>51, <span style="color: #800080;">5</span> <span style="background-color: #38a83d;"> </span></div> <div>60, <span style="color: #800080;">5</span> <span style="background-color: #38a83d;"> </span></div> <div>71, <span style="color: #800080;">3</span> <span style="background-color: #38a83d;"> </span></div> <div>77, <span style="color: #800080;">2</span> <span style="background-color: #38a83d;"> </span></div> <div>...</div>
$\vdots$	

- Problem: List Intersection.
- Instance:  $L_1, \dots, L_n$ .  $n \geq 2$  skip lists of numbers.
- Solution: A sorted list  $R$  of numbers, so that each number occurs in all  $n$  lists.
- Idea:
- (1) Intersection of the two shortest lists  $L_i$  and  $L_j$  to obtain  $R' \supseteq R$ .
  - (2) Iterative intersection of  $R'$  with the remaining lists in ascending order of length.



# Query Processing I

## List Intersection

Algorithm: Intersection of Two Lists.

Input:  $L_1, L_2$ . Skip lists of numbers implemented as singly linked lists.

Output: Sorted list of numbers occurring in both  $L_1$  and  $L_2$ .

*IntersectTwo*( $L_1, L_2$ )

1. Initialization of result list  $R$  and one iterator variable  $x_1$  and  $x_2$  per list.
2. While the iterators point to list entries, process them as follows.
3. If the list entries' keys match, append a merged entry to the result list  $R$ .
4. While the key of  $x_1$  is smaller than that of  $x_2$  advance  $x_1$ .
5. While the key of  $x_2$  is smaller than that of  $x_1$  advance  $x_2$ .
6. Return  $R$ , once an iterator reaches the end of its list.

# Query Processing I

## List Intersection

Algorithm: Intersection of Two Lists.

Input:  $L_1, L_2$ . Skip lists of numbers implemented as singly linked lists.

Output: Sorted list of numbers occurring in both  $L_1$  and  $L_2$ .

*IntersectTwo*( $L_1, L_2$ )

```
1.   $R = list(); x_1 = L_1.head; x_2 = L_2.head$ 
2.  WHILE  $x_1 \neq NIL$  AND  $x_2 \neq NIL$  DO
3.    IF  $x_1.key == x_2.key$  THEN
4.       $R = Insert(R, merge(x_1, x_2))$ 
5.       $x_1 = x_1.next; x_2 = x_2.next$ 
6.    ENDIF
7.    WHILE  $x_1 \neq NIL$  AND  $x_2 \neq NIL$  AND  $x_1.key < x_2.key$  DO
8.      IF CanSkip( $x_1, x_2.key$ ) THEN
9.         $x_1 = Skip(x_1, x_2.key)$ 
10.     ELSE
11.        $x_1 = x_1.next$ 
12.     ENDIF
13.  ENDDO
    :   Like lines 7-13 with  $x_1$  and  $x_2$  exchanged.
21. ENDDO
22. return( $R$ )
```

# Query Processing I

## List Intersection

Algorithm: Intersection of Two Lists.

Input:  $L_1, L_2$ . Skip lists of numbers implemented as singly linked lists.

Output: Sorted list of numbers occurring in both  $L_1$  and  $L_2$ .

*IntersectTwo*( $L_1, L_2$ )

```
1.   $R = list()$ ;  $x_1 = L_1.head$ ;  $x_2 = L_2.head$ 
2.  WHILE  $x_1 \neq NIL$  AND  $x_2 \neq NIL$  DO
3.    IF  $x_1.key == x_2.key$  THEN
4.       $R = Insert(R, merge(x_1, x_2))$ 
5.       $x_1 = x_1.next$ ;  $x_2 = x_2.next$ 
6.    ENDIF
7.     $\vdots$  Like lines 14-20 with  $x_1$  and  $x_2$  exchanged.
14.  WHILE  $x_1 \neq NIL$  AND  $x_2 \neq NIL$  AND  $x_2.key < x_1.key$  DO
15.    IF CanSkip( $x_2, x_1.key$ ) THEN
16.       $x_2 = Skip(x_2, x_1.key)$ 
17.    ELSE
18.       $x_2 = x_2.next$ 
19.    ENDIF
20.  ENDDO
21. ENDDO
22. return( $R$ )
```

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$											
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$											

Execute  $IntersectTwo(L_i, L_j)$ .

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$											
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$											

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = ()$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings																
$\vdots$	$x_i$																
$t_i$	<table><tr><td>2, 4</td><td></td><td></td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td></td><td>19, 7</td><td>23, 5</td><td>28, 6</td><td></td><td>41, 8</td><td>50, 6</td><td>77, 8</td><td></td><td>...</td></tr></table>	2, 4			4, 9	8, 2	16, 1		19, 7	23, 5	28, 6		41, 8	50, 6	77, 8		...
2, 4			4, 9	8, 2	16, 1		19, 7	23, 5	28, 6		41, 8	50, 6	77, 8		...		
$t_j$	<table><tr><td>1, 1</td><td></td><td></td><td>2, 3</td><td>3, 5</td><td>5, 2</td><td></td><td>8, 17</td><td>41, 6</td><td>51, 5</td><td></td><td>60, 5</td><td>71, 3</td><td>77, 2</td><td></td><td>...</td></tr></table>	1, 1			2, 3	3, 5	5, 2		8, 17	41, 6	51, 5		60, 5	71, 3	77, 2		...
1, 1			2, 3	3, 5	5, 2		8, 17	41, 6	51, 5		60, 5	71, 3	77, 2		...		
$\vdots$	$x_j$																

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = ()$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings																
$\vdots$	$x_i$																
$t_i$	<table><tr><td>2, 4</td><td></td><td></td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td></td><td>19, 7</td><td>23, 5</td><td>28, 6</td><td></td><td>41, 8</td><td>50, 6</td><td>77, 8</td><td></td><td>...</td></tr></table>	2, 4			4, 9	8, 2	16, 1		19, 7	23, 5	28, 6		41, 8	50, 6	77, 8		...
2, 4			4, 9	8, 2	16, 1		19, 7	23, 5	28, 6		41, 8	50, 6	77, 8		...		
$t_j$	<table><tr><td>1, 1</td><td></td><td></td><td>2, 3</td><td>3, 5</td><td>5, 2</td><td></td><td>8, 17</td><td>41, 6</td><td>51, 5</td><td></td><td>60, 5</td><td>71, 3</td><td>77, 2</td><td></td><td>...</td></tr></table>	1, 1			2, 3	3, 5	5, 2		8, 17	41, 6	51, 5		60, 5	71, 3	77, 2		...
1, 1			2, 3	3, 5	5, 2		8, 17	41, 6	51, 5		60, 5	71, 3	77, 2		...		
$\vdots$	$x_j$																

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$	$x_i$										
$t_i$	$\boxed{2, \textcolor{violet}{4}}$	$\boxed{4, \textcolor{violet}{9}}$	$\boxed{8, \textcolor{violet}{2}}$	$\boxed{16, \textcolor{violet}{1}}$	$\boxed{19, \textcolor{violet}{7}}$	$\boxed{23, \textcolor{violet}{5}}$	$\boxed{28, \textcolor{violet}{6}}$	$\boxed{41, \textcolor{violet}{8}}$	$\boxed{50, \textcolor{violet}{6}}$	$\boxed{77, \textcolor{violet}{8}}$	...
$t_j$	$\boxed{1, \textcolor{violet}{1}}$	$\boxed{2, \textcolor{violet}{3}}$	$\boxed{3, \textcolor{violet}{5}}$	$\boxed{5, \textcolor{violet}{2}}$	$\boxed{8, \textcolor{violet}{17}}$	$\boxed{41, \textcolor{violet}{6}}$	$\boxed{51, \textcolor{violet}{5}}$	$\boxed{60, \textcolor{violet}{5}}$	$\boxed{71, \textcolor{violet}{3}}$	$\boxed{77, \textcolor{violet}{2}}$	...
$\vdots$	$x_j$										

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots}$



# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$	$x_i$										
$t_i$	$\boxed{2, \textcolor{violet}{4}}$	$\boxed{4, \textcolor{violet}{9}}$	$\boxed{8, \textcolor{violet}{2}}$	$\boxed{16, \textcolor{violet}{1}}$	$\boxed{19, \textcolor{violet}{7}}$	$\boxed{23, \textcolor{violet}{5}}$	$\boxed{28, \textcolor{violet}{6}}$	$\boxed{41, \textcolor{violet}{8}}$	$\boxed{50, \textcolor{violet}{6}}$	$\boxed{77, \textcolor{violet}{8}}$	$\dots$
$t_j$	$\boxed{1, \textcolor{violet}{1}}$	$\boxed{2, \textcolor{violet}{3}}$	$\boxed{3, \textcolor{violet}{5}}$	$\boxed{5, \textcolor{violet}{2}}$	$\boxed{8, \textcolor{violet}{17}}$	$\boxed{41, \textcolor{violet}{6}}$	$\boxed{51, \textcolor{violet}{5}}$	$\boxed{60, \textcolor{violet}{5}}$	$\boxed{71, \textcolor{violet}{3}}$	$\boxed{77, \textcolor{violet}{2}}$	$\dots$
$\vdots$	$x_j$										

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$	$x_i$										
$t_i$	$\boxed{2, \textcolor{violet}{4}}$	$\boxed{4, \textcolor{violet}{9}}$	$\boxed{8, \textcolor{violet}{2}}$	$\boxed{16, \textcolor{violet}{1}}$	$\boxed{19, \textcolor{violet}{7}}$	$\boxed{23, \textcolor{violet}{5}}$	$\boxed{28, \textcolor{violet}{6}}$	$\boxed{41, \textcolor{violet}{8}}$	$\boxed{50, \textcolor{violet}{6}}$	$\boxed{77, \textcolor{violet}{8}}$	$\dots$
$t_j$	$\boxed{1, \textcolor{violet}{1}}$	$\boxed{2, \textcolor{violet}{3}}$	$\boxed{3, \textcolor{violet}{5}}$	$\boxed{5, \textcolor{violet}{2}}$	$\boxed{8, \textcolor{violet}{17}}$	$\boxed{41, \textcolor{violet}{6}}$	$\boxed{51, \textcolor{violet}{5}}$	$\boxed{60, \textcolor{violet}{5}}$	$\boxed{71, \textcolor{violet}{3}}$	$\boxed{77, \textcolor{violet}{2}}$	$\dots$
$\vdots$	$x_j$										

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$	$x_i$										
$t_i$	$\boxed{2, \textcolor{violet}{4}}$	$\boxed{4, \textcolor{violet}{9}}$	$\boxed{8, \textcolor{violet}{2}}$	$\boxed{16, \textcolor{violet}{1}}$	$\boxed{19, \textcolor{violet}{7}}$	$\boxed{23, \textcolor{violet}{5}}$	$\boxed{28, \textcolor{violet}{6}}$	$\boxed{41, \textcolor{violet}{8}}$	$\boxed{50, \textcolor{violet}{6}}$	$\boxed{77, \textcolor{violet}{8}}$	...
$t_j$	$\boxed{1, \textcolor{violet}{1}}$	$\boxed{2, \textcolor{violet}{3}}$	$\boxed{3, \textcolor{violet}{5}}$	$\boxed{5, \textcolor{violet}{2}}$	$\boxed{8, \textcolor{violet}{17}}$	$\boxed{41, \textcolor{violet}{6}}$	$\boxed{51, \textcolor{violet}{5}}$	$\boxed{60, \textcolor{violet}{5}}$	$\boxed{71, \textcolor{violet}{3}}$	$\boxed{77, \textcolor{violet}{2}}$	...
$\vdots$	$x_j$										

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots} \boxed{8, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$	$x_i$										
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$	$x_j$										

Execute *IntersectTwo*( $L_i, L_j$ ).

Result  $R = \boxed{2, \dots} \boxed{8, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$	$x_i$										
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$	$x_j$										

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots} \boxed{8, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$											
	$x_i$										
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$	$x_j$										

Execute *IntersectTwo*( $L_i, L_j$ ).

Result  $R = \boxed{2, \dots} \boxed{8, \dots} \boxed{41, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$											
	$x_i$										
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$	$x_j$										

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots} \boxed{8, \dots} \boxed{41, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$											
											$x_i$
$t_i$	$2, \textcolor{violet}{4}$	$4, \textcolor{violet}{9}$	$8, \textcolor{violet}{2}$	$16, \textcolor{violet}{1}$	$19, \textcolor{violet}{7}$	$23, \textcolor{violet}{5}$	$28, \textcolor{violet}{6}$	$41, \textcolor{violet}{8}$	$50, \textcolor{violet}{6}$	$77, \textcolor{violet}{8}$	...
$t_j$	$1, \textcolor{violet}{1}$	$2, \textcolor{violet}{3}$	$3, \textcolor{violet}{5}$	$5, \textcolor{violet}{2}$	$8, \textcolor{violet}{17}$	$41, \textcolor{violet}{6}$	$51, \textcolor{violet}{5}$	$60, \textcolor{violet}{5}$	$71, \textcolor{violet}{3}$	$77, \textcolor{violet}{2}$	...
$\vdots$											
											$x_j$

Execute *IntersectTwo*( $L_i, L_j$ ).

Result  $R = \boxed{2, \dots} \boxed{8, \dots} \boxed{41, \dots}$



# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings										
$\vdots$											
	$x_i$										
$t_i$	$\boxed{2, \textcolor{violet}{4}}$	$\boxed{4, \textcolor{violet}{9}}$	$\boxed{8, \textcolor{violet}{2}}$	$\boxed{16, \textcolor{violet}{1}}$	$\boxed{19, \textcolor{violet}{7}}$	$\boxed{23, \textcolor{violet}{5}}$	$\boxed{28, \textcolor{violet}{6}}$	$\boxed{41, \textcolor{violet}{8}}$	$\boxed{50, \textcolor{violet}{6}}$	$\boxed{77, \textcolor{violet}{8}}$	...
$t_j$	$\boxed{1, \textcolor{violet}{1}}$	$\boxed{2, \textcolor{violet}{3}}$	$\boxed{3, \textcolor{violet}{5}}$	$\boxed{5, \textcolor{violet}{2}}$	$\boxed{8, \textcolor{violet}{17}}$	$\boxed{41, \textcolor{violet}{6}}$	$\boxed{51, \textcolor{violet}{5}}$	$\boxed{60, \textcolor{violet}{5}}$	$\boxed{71, \textcolor{violet}{3}}$	$\boxed{77, \textcolor{violet}{2}}$	...
$\vdots$											
	$x_j$										

Execute *IntersectTwo*( $L_i, L_j$ ).

Result  $R = \boxed{2, \dots} \boxed{8, \dots} \boxed{41, \dots} \boxed{77, \dots}$

# Query Processing I

## List Intersection: Example

Given an index with postings  $\boxed{k, \textcolor{violet}{tf}(t, d_k)}$ , two postlists  $L_i, L_j$  for terms  $t_i, t_j$ , and the query  $q = t_i \wedge t_j$ :

$T$	Postings											
$\vdots$												$x_i$
$t_i$	$\boxed{2, \textcolor{violet}{4}}$	$\boxed{4, \textcolor{violet}{9}}$	$\boxed{8, \textcolor{violet}{2}}$	$\boxed{16, \textcolor{violet}{1}}$	$\boxed{19, \textcolor{violet}{7}}$	$\boxed{23, \textcolor{violet}{5}}$	$\boxed{28, \textcolor{violet}{6}}$	$\boxed{41, \textcolor{violet}{8}}$	$\boxed{50, \textcolor{violet}{6}}$	$\boxed{77, \textcolor{violet}{8}}$	$\dots$	
$t_j$	$\boxed{1, \textcolor{violet}{1}}$	$\boxed{2, \textcolor{violet}{3}}$	$\boxed{3, \textcolor{violet}{5}}$	$\boxed{5, \textcolor{violet}{2}}$	$\boxed{8, \textcolor{violet}{17}}$	$\boxed{41, \textcolor{violet}{6}}$	$\boxed{51, \textcolor{violet}{5}}$	$\boxed{60, \textcolor{violet}{5}}$	$\boxed{71, \textcolor{violet}{3}}$	$\boxed{77, \textcolor{violet}{2}}$	$\dots$	
$\vdots$												$x_j$

Execute  $IntersectTwo(L_i, L_j)$ .

Result  $R = \boxed{2, \dots} \boxed{8, \dots} \boxed{41, \dots} \boxed{77, \dots}$

## Remarks:

- ❑ Postlists are usually too large to fit in main memory, so iterating them brings performance benefits.
- ❑ The *key* attribute stores the document identifier of a posting.
- ❑ The *merge* function returns a posting merged from the two postings passed in. It merges the potentially stored term weights and other information stored in them.
- ❑ The *next* attribute stores the successive posting.
- ❑ The *CanSkip* function checks whether the current posting contains skip information and whether a target with a document identifier less than or equal to the *key* value passed is available.
- ❑ The *Skip* function returns the posting that is closest to but less than or equal to the *key* value passed.

# Query Processing I

## List Intersection

Algorithm: Intersect Many Lists.

Input:  $L_1, \dots, L_n$ . Skip lists of numbers implemented as singly linked lists.

Output: Sorted list of numbers occurring in all  $L_1, \dots, L_n$ .

*IntersectMany*( $L_1, \dots, L_n$ )

// Sort by list length.

1.  $H = \text{BuildMinHeap}(L_1, \dots, L_n);$

2.  $R = \text{ExtractMin}(H)$

3. **WHILE**  $|H| > 0$  **DO**

4.      $L_{\min} = \text{ExtractMin}(H)$

5.      $R = \text{IntersectTwo}(R, L_{\min})$

6. **ENDDO**

7. *return*( $R$ )

Why are lists intersected in ascending order of list length?

# Query Processing I

## List Intersection

Algorithm: Intersect Many Lists.

Input:  $L_1, \dots, L_n$ . Skip lists of numbers implemented as singly linked lists.

Output: Sorted list of numbers occurring in all  $L_1, \dots, L_n$ .

*IntersectMany*( $L_1, \dots, L_n$ )

```
// Sort by list length.  
1.  $H = \text{BuildMinHeap}(L_1, \dots, L_n);$   
2.  $R = \text{ExtractMin}(H)$   
3. WHILE  $|H| > 0$  DO  
4.    $L_{\min} = \text{ExtractMin}(H)$   
5.    $R = \text{IntersectTwo}(R, L_{\min})$   
6. ENDDO  
7. return( $R$ )
```

Observations:

- ❑ The amount of memory required to store the result list  $R$  is bounded by the shortest list from  $L_1, \dots, L_n$ .
- ❑ The smaller the result list  $R$ , the more effective are the skip pointers.
- ❑ Hard disk seeking is minimized since every list is read sequentially.

# Query Processing I

## Proximity Queries

Given a query  $q = t_i / \epsilon t_j$ , retrieve documents in which  $t_i$  and  $t_j$  are in close proximity, i.e., within an  $\epsilon$ -environment of one another, where  $\epsilon \geq 1$  terms.

# Query Processing I

## Proximity Queries

Given a query  $q = t_i / \epsilon t_j$ , retrieve documents in which  $t_i$  and  $t_j$  are in close **proximity**, i.e., within an  **$\epsilon$ -environment** of one another, where  $\epsilon \geq 1$  terms.

Processing proximity queries requires term positions in postings:

`<document>   [<weights>]   [<positions>]   [...]`

# Query Processing I

## Proximity Queries

Given a query  $q = t_i / \epsilon t_j$ , retrieve documents in which  $t_i$  and  $t_j$  are in close **proximity**, i.e., within an  **$\epsilon$ -environment** of one another, where  $\epsilon \geq 1$  terms.

Processing proximity queries requires term positions in postings:

`<document> [ <weights> ] [ <positions> ] [ ... ]`

Example:

$d =$  “You cannot end a sentence with because because because is a conjunction.”  
          1      2      3      4      5      6      7      8      9     10 11     12

Posting for “because” and  $d$ :

$d,$  3, (7, 8, 9)

Posting for “sentence” and  $d$ :

$d,$  1, (5)

In  $d$ , “because” is in a 2-environment of {“sentence”, “with”, “because”, “is”, “a”}.



# Query Processing I

## Proximity Queries

Algorithm: Position List Intersection.

Input:  $A_1, A_2$ . Sorted arrays of positions of two terms  $t_1, t_2$  in a document  $d$ .  
 $\epsilon$ . Maximal term distance.

Output: For each position in  $A_1$ , the positions from  $A_2$  within an  $\epsilon$ -environment.

*IntersectPositions*( $A_1, A_2, \epsilon$ )

```
1.  $R = \text{map}()$ 
2. FOR  $i = 1$  TO  $A_1.\text{length}$  DO
3.    $R' = \text{list}()$ 
4.   FOR  $j = 1$  TO  $A_2.\text{length}$  DO
5.     IF  $|A_1[i] - A_2[j]| \leq \epsilon$  THEN
6.        $\text{insert}(R', A_2[j])$ 
7.     ELSE IF  $A_2[j] > A_1[i]$  THEN
8.        $\text{break}$ 
9.     ENDIF
10.  ENDDO
11.   $\text{insert}(R, A_1[i], R')$ 
12. ENDDO
13.  $\text{return}(R)$ 
```

Remarks:

- ❑ Pruning unnecessary comparisons  
Lines 7–9: Stop comparing once the  $j$ -th position in  $A_2$  exceeds the  $i$ -th position in  $A_1$  by more than  $\epsilon$ . The difference can never get smaller than  $\epsilon$  again.
- ❑ Integration into postlist intersection  
The if-statement of Line 3 of [\*IntersectTwo\*](#) then additionally checks whether *IntersectPositions* returns a non-empty result.

# Query Processing I

## Phrase Queries

Given a phrase query  $q = "t_1 \dots t_m"$ , retrieve documents in which the terms  $t_1, \dots, t_m$  occur in the same order as in the query  $q$ .

# Query Processing I

## Phrase Queries

Given a phrase query  $q = "t_1 \dots t_m"$ , retrieve documents in which the terms  $t_1, \dots, t_m$  occur in the same order as in the query  $q$ .

Processing phrase queries requires term positions in postings.

Example:

$T$	Postings
to	... 4, 250, (... 133, 137, ...) ...
be	... 4, 125, (... 134, 138, ...) ...
or	... 4, 40, (... 135, ...) ...
not	... 4, 15, (... 136, ...) ...

What phrase does document 4 contain?

# Query Processing I

## Phrase Queries

Given a phrase query  $q = "t_1 \dots t_m"$ , retrieve documents in which the terms  $t_1, \dots, t_m$  occur in the same order as in the query  $q$ .

Processing phrase queries requires term positions in postings.

Example:

$T$	Postings
to	... <span style="border: 1px solid black; padding: 2px;">4, 250, (... , 133, 137, ...)</span> ...
be	... <span style="border: 1px solid black; padding: 2px;">4, 125, (... , 134, 138, ...)</span> ...
or	... <span style="border: 1px solid black; padding: 2px;">4, 40, (... , 135, ...)</span> ...
not	... <span style="border: 1px solid black; padding: 2px;">4, 15, (... , 136, ...)</span> ...

Document 4 contains the phrase  
to be or not to be  
at term positions 133–138.

Observations:

- ❑ Processing phrase queries can be reduced to the list intersection problem.  
*Algorithms `IntersectMany` and `IntersectTwo` can be adjusted to process phrase queries.*
- ❑ The additional run time for phrase processing is in  $O(\sum_{d \in \text{IntersectMany}(L_t: t \in q)} |d|)$ .

# Query Processing I

## Phrase Queries

Given a phrase query  $q = "t_1 \dots t_m"$ , retrieve documents in which the terms  $t_1, \dots, t_m$  occur in the same order as in the query  $q$ .

To speed up phrase search, ***n*-grams** can be used as index terms.

Example:

$T$	Postings	
to be	... <table><tr><td>4, 80, (... , 133, 137, ...)</td></tr></table> ...	4, 80, (... , 133, 137, ...)
4, 80, (... , 133, 137, ...)		
be or	... <table><tr><td>4, 55, (... , 134, ...)</td></tr></table> ...	4, 55, (... , 134, ...)
4, 55, (... , 134, ...)		
or not	... <table><tr><td>4, 20, (... , 135, ...)</td></tr></table> ...	4, 20, (... , 135, ...)
4, 20, (... , 135, ...)		
not to	... <table><tr><td>4, 7, (... , 136, ...)</td></tr></table> ...	4, 7, (... , 136, ...)
4, 7, (... , 136, ...)		

Document 4 contains the phrase  
to be or not to be  
at term positions 133–138.

How much faster can phrase queries be processed?

# Query Processing I

## Phrase Queries

Given a phrase query  $q = "t_1 \dots t_m"$ , retrieve documents in which the terms  $t_1, \dots, t_m$  occur in the same order as in the query  $q$ .

To speed up phrase search,  **$n$ -grams** can be used as index terms.

Example:

$T$	Postings	
to be	... <table><tr><td>4, 80, (... , 133, 137, ...)</td></tr></table> ...	4, 80, (... , 133, 137, ...)
4, 80, (... , 133, 137, ...)		
be or	... <table><tr><td>4, 55, (... , 134, ...)</td></tr></table> ...	4, 55, (... , 134, ...)
4, 55, (... , 134, ...)		
or not	... <table><tr><td>4, 20, (... , 135, ...)</td></tr></table> ...	4, 20, (... , 135, ...)
4, 20, (... , 135, ...)		
not to	... <table><tr><td>4, 7, (... , 136, ...)</td></tr></table> ...	4, 7, (... , 136, ...)
4, 7, (... , 136, ...)		

Document 4 contains the phrase  
to be or not to be  
at term positions 133–138.

Observations:

- ❑ The time to process phrase queries of length at least  $n$  is divided by  $n$ .  
Only non-overlapping  $n$ -grams need to be intersected.
- ❑ Maintaining an index with  $n$ -grams and/or common phrases as index terms speeds up non-phrase queries as well.

## Remarks:

- ❑ The space requirements of a positional index are 2–4 times that of a nonpositional index.
- ❑ Most basic retrieval models do not directly employ positional information. If keyword proximity is a desired feature in a retrieval system using a basic retrieval model, positional information usually is implemented as an additional relevance signal or as a prior probability for a document.

# Chapter IR:II

## II. Indexing

- ❑ Indexing Basics
- ❑ Inverted Index
- ❑ Query Processing I
- ❑ Query Processing II
- ❑ Index Construction
- ❑ Index Compression
- ❑ Size Estimation



# Query Processing II

## Retrieval Types

Query processing can be based on two basic approaches:

- ❑ **Set retrieval**

A query induces a subset of the indexed documents which is considered relevant.  
Important applications: e-discovery, patent search, systematic reviews.

- ❑ **Ranked retrieval**

A query induces a ranking among all indexed documents in descending order of relevance.

Ranked retrieval is the norm in virtually all modern search engines.

# Query Processing II

## Relevance Scoring (Recap)

Quantification of the relevance of an indexed document  $d$  to a query  $q$ .

# Query Processing II

## Relevance Scoring (Recap)

Quantification of the relevance of an indexed document  $d$  to a query  $q$ .

Let  $t \in T$  denote a term  $t$  from the terminology  $T$  of index terms, and let  $\omega_X : T \times X \rightarrow \mathbf{R}$  denote a term weighting function, where  $X$  may be a set of documents  $D$  or a set of queries  $Q$ . Then the most basic relevance function  $\rho$  is:

$$\rho(q, d) = \sum_{t \in T} \omega_Q(t, q) \cdot \omega_D(t, d),$$

where  $\omega_Q(t, q)$  and  $\omega_D(t, d)$  are term weights indicating the importance of  $t$  for the query  $q \in Q$  and the document  $d \in D$ , respectively.

# Query Processing II

## Relevance Scoring (Recap)

Quantification of the relevance of an indexed document  $d$  to a query  $q$ .

Let  $t \in T$  denote a term  $t$  from the terminology  $T$  of index terms, and let  $\omega_X : T \times X \rightarrow \mathbf{R}$  denote a term weighting function, where  $X$  may be a set of documents  $D$  or a set of queries  $Q$ . Then the most basic relevance function  $\rho$  is:

$$\rho(q, d) = \sum_{t \in T} \omega_Q(t, q) \cdot \omega_D(t, d),$$

where  $\omega_Q(t, q)$  and  $\omega_D(t, d)$  are term weights indicating the importance of  $t$  for the query  $q \in Q$  and the document  $d \in D$ , respectively.

Observations:

- ❑ A term  $t$  may have importance, and hence non-zero weights, for a query  $q$  or document  $d$  despite not occurring in them. Example: synonyms.
- ❑ The majority of terms from  $T$  will have insignificant importance to both.
- ❑ The term weights  $\omega_D(t, d)$  can be pre-computed and indexed.
- ❑ The term weights  $\omega_Q(t, q)$  must be computed on the fly.

# Query Processing II

## Query Semantics for Ranked Retrieval

Keyword queries have Boolean semantics that is either implicitly specified by user behavior and expectations or explicitly specified.

We distinguish four types:

- ❑ Single-term queries
- ❑ Disjunctive multi-term queries  
Only Boolean OR connectives. Example: Antony  $\vee$  Brutus  $\vee$  Calpurnia.
- ❑ Conjunctive multi-term queries  
Only Boolean AND connectives. Example: Antony  $\wedge$  Brutus  $\wedge$  Calpurnia.
  - + Constraint: Proximity  
Example: Antony  $/\epsilon$  Caesar
  - + Constraint: Phrase  
Example: “Antony and Caesar”
- ❑ “Complex” Boolean multi-term queries  
Remainder of Boolean formulas. Example: (Antony  $\vee$  Caesar)  $\wedge \neg$  Calpurnia.  
Can be normalized to disjunctive or conjunctive normal form.

# Query Processing II

## Single-Term Queries

Given a single-term query  $q = t$ , the optimal postlist ordering is by term weight.

Example:

$T$	Postings (ordered by document identifier)										
$\vdots$											
$t_i$	2, 4	4, 9	8, 2	16, 1	19, 7	23, 5	28, 6	41, 8	50, 6	77, 8	...
$t_j$	1, 1	2, 3	3, 5	5, 2	8, 17	41, 6	51, 5	60, 5	71, 3	77, 2	...
$\vdots$											

Worst case: The last document of the postlist is the most relevant one.  
The whole postlist must be examined.

# Query Processing II

## Single-Term Queries

Given a single-term query  $q = t$ , the optimal postlist ordering is by term weight.

Example:

<i>T</i>	Postings (ordered by term weight)										
⋮											
$t_i$	4, 9	41, 8	77, 8	19, 7	28, 6	50, 6	23, 5	2, 4	8, 2	16, 1	...
$t_j$	8, 17	41, 6	3, 5	51, 5	60, 5	2, 3	71, 3	5, 2	77, 2	1, 1	...
⋮											

Best case: The document whose content is best represented by the term  $t$  is the one with the highest term weight. A partial examination of the postlist suffices.

Including a skip list in a postlist ordered by term weights may not be useful.

# Query Processing II

## Disjunctive Queries

In general, a query  $q$  is processed as a **disjunctive query**, where each term  $t_i \in q$  may or may not occur in a relevant document  $d$ , as long as at least one  $t_i$  occurs.

### Document-at-a-time scoring

- ❑ Precondition: a total order of documents in the index's postlists is enforced  
Ordering criterion: document ID or document quality
- ❑ Parallel traversal of query term postlists, document ID by document ID.
- ❑ Each document's score is instantly complete, but the ranking only at the end.
- ❑ Concurrent disk IO overhead increases with query length.



# Query Processing II

## Disjunctive Queries

In general, a query  $q$  is processed as a **disjunctive query**, where each term  $t_i \in q$  may or may not occur in a relevant document  $d$ , as long as at least one  $t_i$  occurs.

### Document-at-a-time scoring

- ❑ Precondition: a total order of documents in the index's postlists is enforced  
Ordering criterion: document ID or document quality
- ❑ Parallel traversal of query term postlists, document ID by document ID.
- ❑ Each document's score is instantly complete, but the ranking only at the end.
- ❑ Concurrent disk IO overhead increases with query length.

### Term-at-a-time scoring

- ❑ Iterative traversal of query term postlists (e.g., in order of term frequency).
- ❑ Temporary query postlist contains candidate documents.
- ❑ As document scores accumulate, an approximate ranking becomes available.
- ❑ More main memory required for maintaining temporary postlist.

Safe and unsafe optimizations exist (e.g., to stop the search early).

## Remarks:

- ❑ Web search engines often return results without some of a query's terms for very specific queries, indicating a disjunctive interpretation. Nevertheless, many retrieval models assign higher scores to documents matching more of a query's terms, leaning toward a “conjunctive” interpretation at least for the (visible) top results.

# Query Processing II

## Disjunctive Queries

Algorithm: Document-at-a-time Scoring.

Input:  $L_1, \dots, L_m$ . The postlists of the terms  $t_1, \dots, t_m$  of query  $q$ .  
 $\mathbf{q}$ . Representation of query  $q$ , e.g., as array of  $m$  term weights.

Output: A list of documents in  $D$ , sorted in descending order of relevance to  $q$ .

*DAATScoring*( $L_1, \dots, L_m, \mathbf{q}$ )

1. Initialization of result list  $R$  as priority queue, and postlist iterator variables.
2. While not all postlists have been processed, repeat the following steps.
3. Determine the smallest document identifier  $d$  to which the iterators point.
4. Collect all term weights of  $d$  in an array  $\mathbf{d}$ .
5. Calculate the relevance score  $\rho(\mathbf{q}, \mathbf{d})$  and insert it in  $R$ .
6. Advance all iterators pointing to  $d$ .
7. Return the list of scored documents  $R$ .

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```
1.  $R = PriorityQueue()$ 
2.  $x_1 = L_1.head; \dots; x_m = L_m.head$ 
3.  $continue = TRUE$ 
4. WHILE  $continue$  DO
5.    $d = \min_{i \in [1, m]}(x_i.key)$ 
6.    $\mathbf{d} = Array(|q|)$ 
7.   FOR  $i \in [1, m]$  DO
8.     IF  $x_i \neq NIL$  AND  $x_i.key = d$  THEN
9.        $\mathbf{d}[i] = x_i.weight$ 
10.    ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $Insert(R, record(d, r))$ 
14.   $continue = FALSE$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq NIL$  AND  $x_i.key = d$  THEN
17.       $x_i = x_i.next$ 
18.    ENDIF
19.    IF  $x_i \neq NIL$  THEN
20.       $continue = TRUE$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $return(R)$ 
```

Example:

$T$	Postings						
$\vdots$							
$t_i$	<table><tr><td>1, 4</td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td>19, 7</td><td>...</td></tr></table>	1, 4	4, 9	8, 2	16, 1	19, 7	...
1, 4	4, 9	8, 2	16, 1	19, 7	...		
$\vdots$							
$t_j$	<table><tr><td>1, 1</td><td>2, 3</td><td>5, 5</td><td>7, 2</td><td>8, 8</td><td>...</td></tr></table>	1, 1	2, 3	5, 5	7, 2	8, 8	...
1, 1	2, 3	5, 5	7, 2	8, 8	...		
$\vdots$							
$t_k$	<table><tr><td>1, 2</td><td>2, 4</td><td>5, 1</td><td>6, 3</td><td>8, 5</td><td>...</td></tr></table>	1, 2	2, 4	5, 1	6, 3	8, 5	...
1, 2	2, 4	5, 1	6, 3	8, 5	...		
$\vdots$							

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} =$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```
1.  $R = \text{PriorityQueue}()$ 
2.  $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.  $\text{continue} = \text{TRUE}$ 
4. WHILE  $\text{continue}$  DO
5.    $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.    $\mathbf{d} = \text{Array}(|q|)$ 
7.   FOR  $i \in [1, m]$  DO
8.     IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.        $\mathbf{d}[i] = x_i.\text{weight}$ 
10.    ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 
```

Example:

$T$	Postings						
$\vdots$	$x_i$						
$t_i$	<table><tr><td>1, 4</td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td>19, 7</td><td>...</td></tr></table>	1, 4	4, 9	8, 2	16, 1	19, 7	...
1, 4	4, 9	8, 2	16, 1	19, 7	...		
$\vdots$	$x_j$						
$t_j$	<table><tr><td>1, 1</td><td>2, 3</td><td>5, 5</td><td>7, 2</td><td>8, 8</td><td>...</td></tr></table>	1, 1	2, 3	5, 5	7, 2	8, 8	...
1, 1	2, 3	5, 5	7, 2	8, 8	...		
$\vdots$	$x_k$						
$t_k$	<table><tr><td>1, 2</td><td>2, 4</td><td>5, 1</td><td>6, 3</td><td>8, 5</td><td>...</td></tr></table>	1, 2	2, 4	5, 1	6, 3	8, 5	...
1, 2	2, 4	5, 1	6, 3	8, 5	...		
$\vdots$							

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 4 \\ \vdots \\ 1 \\ \vdots \\ 2 \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.   $R = \text{PriorityQueue}()$ 
2.   $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.   $\text{continue} = \text{TRUE}$ 
4.  WHILE  $\text{continue}$  DO
5.     $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.     $\mathbf{d} = \text{Array}(|q|)$ 
7.    FOR  $i \in [1, m]$  DO
8.      IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.         $\mathbf{d}[i] = x_i.\text{weight}$ 
10.     ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings
$\vdots$	$x_i$
$t_i$	<span>1, 4</span> <span>4, 9</span> <span>8, 2</span> <span>16, 1</span> <span>19, 7</span> ...
$\vdots$	$x_j$
$t_j$	<span>1, 1</span> <span>2, 3</span> <span>5, 5</span> <span>7, 2</span> <span>8, 8</span> ...
$\vdots$	$x_k$
$t_k$	<span>1, 2</span> <span>2, 4</span> <span>5, 1</span> <span>6, 3</span> <span>8, 5</span> ...
$\vdots$	

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 0 \\ \vdots \\ 3 \\ \vdots \\ 4 \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.  $R = \text{PriorityQueue}()$ 
2.  $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.  $\text{continue} = \text{TRUE}$ 
4. WHILE  $\text{continue}$  DO
5.    $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.    $\mathbf{d} = \text{Array}(|q|)$ 
7.   FOR  $i \in [1, m]$  DO
8.     IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.        $\mathbf{d}[i] = x_i.\text{weight}$ 
10.    ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings						
$\vdots$	$x_i$						
$t_i$	<table><tr><td>1, 4</td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td>19, 7</td><td>...</td></tr></table>	1, 4	4, 9	8, 2	16, 1	19, 7	...
1, 4	4, 9	8, 2	16, 1	19, 7	...		
$\vdots$	$x_j$						
$t_j$	<table><tr><td>1, 1</td><td>2, 3</td><td>5, 5</td><td>7, 2</td><td>8, 8</td><td>...</td></tr></table>	1, 1	2, 3	5, 5	7, 2	8, 8	...
1, 1	2, 3	5, 5	7, 2	8, 8	...		
$\vdots$	$x_k$						
$t_k$	<table><tr><td>1, 2</td><td>2, 4</td><td>5, 1</td><td>6, 3</td><td>8, 5</td><td>...</td></tr></table>	1, 2	2, 4	5, 1	6, 3	8, 5	...
1, 2	2, 4	5, 1	6, 3	8, 5	...		
$\vdots$							

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 9 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```
1.  $R = \text{PriorityQueue}()$ 
2.  $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.  $\text{continue} = \text{TRUE}$ 
4. WHILE  $\text{continue}$  DO
5.    $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.    $\mathbf{d} = \text{Array}(|q|)$ 
7.   FOR  $i \in [1, m]$  DO
8.     IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.        $\mathbf{d}[i] = x_i.\text{weight}$ 
10.    ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 
```

Example:

$T$	Postings						
$\vdots$	$x_i$						
$t_i$	<table><tr><td>1, 4</td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td>19, 7</td><td>...</td></tr></table>	1, 4	4, 9	8, 2	16, 1	19, 7	...
1, 4	4, 9	8, 2	16, 1	19, 7	...		
$\vdots$	$x_j$						
$t_j$	<table><tr><td>1, 1</td><td>2, 3</td><td>5, 5</td><td>7, 2</td><td>8, 8</td><td>...</td></tr></table>	1, 1	2, 3	5, 5	7, 2	8, 8	...
1, 1	2, 3	5, 5	7, 2	8, 8	...		
$\vdots$	$x_k$						
$t_k$	<table><tr><td>1, 2</td><td>2, 4</td><td>5, 1</td><td>6, 3</td><td>8, 5</td><td>...</td></tr></table>	1, 2	2, 4	5, 1	6, 3	8, 5	...
1, 2	2, 4	5, 1	6, 3	8, 5	...		
$\vdots$							

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 0 \\ \vdots \\ 5 \\ \vdots \\ 1 \\ \vdots \end{pmatrix}$$



# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.   $R = \text{PriorityQueue}()$ 
2.   $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.   $\text{continue} = \text{TRUE}$ 
4.  WHILE  $\text{continue}$  DO
5.     $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.     $\mathbf{d} = \text{Array}(|q|)$ 
7.    FOR  $i \in [1, m]$  DO
8.      IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.         $\mathbf{d}[i] = x_i.\text{weight}$ 
10.     ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings						
$\vdots$	$x_i$						
$t_i$	<table><tr><td>1, 4</td><td>4, 9</td><td>8, 2</td><td>16, 1</td><td>19, 7</td><td>...</td></tr></table>	1, 4	4, 9	8, 2	16, 1	19, 7	...
1, 4	4, 9	8, 2	16, 1	19, 7	...		
$\vdots$	$x_j$						
$t_j$	<table><tr><td>1, 1</td><td>2, 3</td><td>5, 5</td><td>7, 2</td><td>8, 8</td><td>...</td></tr></table>	1, 1	2, 3	5, 5	7, 2	8, 8	...
1, 1	2, 3	5, 5	7, 2	8, 8	...		
$\vdots$	$x_k$						
$t_k$	<table><tr><td>1, 2</td><td>2, 4</td><td>5, 1</td><td>6, 3</td><td>8, 5</td><td>...</td></tr></table>	1, 2	2, 4	5, 1	6, 3	8, 5	...
1, 2	2, 4	5, 1	6, 3	8, 5	...		
$\vdots$							

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 3 \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.  $R = PriorityQueue()$ 
2.  $x_1 = L_1.head; \dots; x_m = L_m.head$ 
3.  $continue = TRUE$ 
4. WHILE  $continue$  DO
5.    $d = \min_{i \in [1, m]}(x_i.key)$ 
6.    $\mathbf{d} = Array(|q|)$ 
7.   FOR  $i \in [1, m]$  DO
8.     IF  $x_i \neq NIL$  AND  $x_i.key = d$  THEN
9.        $\mathbf{d}[i] = x_i.weight$ 
10.    ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $Insert(R, record(d, r))$ 
14.   $continue = FALSE$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq NIL$  AND  $x_i.key = d$  THEN
17.       $x_i = x_i.next$ 
18.    ENDIF
19.    IF  $x_i \neq NIL$  THEN
20.       $continue = TRUE$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $return(R)$ 

```

Example:

$T$	Postings
$\vdots$	$x_i$
$t_i$	<span>1, 4</span> <span>4, 9</span> <span>8, 2</span> <span>16, 1</span> <span>19, 7</span> ...
$\vdots$	$x_j$
$t_j$	<span>1, 1</span> <span>2, 3</span> <span>5, 5</span> <span>7, 2</span> <span>8, 8</span> ...
$\vdots$	$x_k$
$t_k$	<span>1, 2</span> <span>2, 4</span> <span>5, 1</span> <span>6, 3</span> <span>8, 5</span> ...
$\vdots$	

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 0 \\ \vdots \\ 2 \\ \vdots \\ 0 \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.   $R = \text{PriorityQueue}()$ 
2.   $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.   $\text{continue} = \text{TRUE}$ 
4.  WHILE  $\text{continue}$  DO
5.     $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.     $\mathbf{d} = \text{Array}(|q|)$ 
7.    FOR  $i \in [1, m]$  DO
8.      IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.         $\mathbf{d}[i] = x_i.\text{weight}$ 
10.     ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings
$\vdots$	$x_i$
$t_i$	<span>1, 4</span> <span>4, 9</span> <span>8, 2</span> <span>16, 1</span> <span>19, 7</span> ...
$\vdots$	$x_j$
$t_j$	<span>1, 1</span> <span>2, 3</span> <span>5, 5</span> <span>7, 2</span> <span>8, 8</span> ...
$\vdots$	$x_k$
$t_k$	<span>1, 2</span> <span>2, 4</span> <span>5, 1</span> <span>6, 3</span> <span>8, 5</span> ...
$\vdots$	

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 2 \\ \vdots \\ 8 \\ \vdots \\ 5 \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.   $R = \text{PriorityQueue}()$ 
2.   $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.   $\text{continue} = \text{TRUE}$ 
4.  WHILE  $\text{continue}$  DO
5.     $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.     $\mathbf{d} = \text{Array}(|q|)$ 
7.    FOR  $i \in [1, m]$  DO
8.      IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.         $\mathbf{d}[i] = x_i.\text{weight}$ 
10.     ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings
$\vdots$	$x_i$
$t_i$	<span>1, 4</span> <span>4, 9</span> <span>8, 2</span> <span>16, 1</span> <span>19, 7</span> ...
$\vdots$	$x_j$
$t_j$	<span>1, 1</span> <span>2, 3</span> <span>5, 5</span> <span>7, 2</span> <span>8, 8</span> ...
$\vdots$	$x_k$
$t_k$	<span>1, 2</span> <span>2, 4</span> <span>5, 1</span> <span>6, 3</span> <span>8, 5</span> ...
$\vdots$	

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 1 \\ \vdots \\ w_j \\ \vdots \\ w_k \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.   $R = \text{PriorityQueue}()$ 
2.   $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.   $\text{continue} = \text{TRUE}$ 
4.  WHILE  $\text{continue}$  DO
5.     $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.     $\mathbf{d} = \text{Array}(|q|)$ 
7.    FOR  $i \in [1, m]$  DO
8.      IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.         $\mathbf{d}[i] = x_i.\text{weight}$ 
10.     ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings
$\vdots$	$x_i$
$t_i$	<span>1, 4</span> <span>4, 9</span> <span>8, 2</span> <span>16, 1</span> <span>19, 7</span> ...
$\vdots$	$x_j$
$t_j$	<span>1, 1</span> <span>2, 3</span> <span>5, 5</span> <span>7, 2</span> <span>8, 8</span> ...
$\vdots$	$x_k$
$t_k$	<span>1, 2</span> <span>2, 4</span> <span>5, 1</span> <span>6, 3</span> <span>8, 5</span> ...
$\vdots$	

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ 7 \\ \vdots \\ w_j \\ \vdots \\ w_k \\ \vdots \end{pmatrix}$$

# Query Processing II

## Disjunctive Queries

$DAATScoring(L_1, \dots, L_m, \mathbf{q})$

```

1.   $R = \text{PriorityQueue}()$ 
2.   $x_1 = L_1.\text{head}; \dots; x_m = L_m.\text{head}$ 
3.   $\text{continue} = \text{TRUE}$ 
4.  WHILE  $\text{continue}$  DO
5.     $d = \min_{i \in [1, m]}(x_i.\text{key})$ 
6.     $\mathbf{d} = \text{Array}(|q|)$ 
7.    FOR  $i \in [1, m]$  DO
8.      IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
9.         $\mathbf{d}[i] = x_i.\text{weight}$ 
10.     ENDIF
11.  ENDDO
12.   $r = \rho(\mathbf{q}, \mathbf{d})$ 
13.   $\text{Insert}(R, \text{record}(d, r))$ 
14.   $\text{continue} = \text{FALSE}$ 
15.  FOR  $i \in [1, m]$  DO
16.    IF  $x_i \neq \text{NIL}$  AND  $x_i.\text{key} = d$  THEN
17.       $x_i = x_i.\text{next}$ 
18.    ENDIF
19.    IF  $x_i \neq \text{NIL}$  THEN
20.       $\text{continue} = \text{TRUE}$ 
21.    ENDIF
22.  ENDDO
23. ENDDO
24.  $\text{return}(R)$ 

```

Example:

$T$	Postings
$\vdots$	$x_i$
$t_i$	<span>1, 4</span> <span>4, 9</span> <span>8, 2</span> <span>16, 1</span> <span>19, 7</span> ...
$\vdots$	$x_j$
$t_j$	<span>1, 1</span> <span>2, 3</span> <span>5, 5</span> <span>7, 2</span> <span>8, 8</span> ...
$\vdots$	$x_k$
$t_k$	<span>1, 2</span> <span>2, 4</span> <span>5, 1</span> <span>6, 3</span> <span>8, 5</span> ...
$\vdots$	

$$\mathbf{q} = \begin{pmatrix} \vdots \\ 5 \\ \vdots \\ 8 \\ \vdots \\ 3 \\ \vdots \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} \vdots \\ w_i \\ \vdots \\ w_j \\ \vdots \\ w_k \\ \vdots \end{pmatrix}$$

## Remarks:

- ❑ DAAT = Document at a time
- ❑ We distinguish between a real-world query  $q$  and its computer representation  $\mathbf{q}$ . Likewise, document (identifier)  $d$ 's representation is  $\mathbf{d}$ . More complex representations can be imagined than the array-of-weights representations exemplified.
- ❑ Relevance function  $\rho(\mathbf{q}, \mathbf{d})$  maps pairs of document and query representations to a real-valued score indicating document  $d$ 's relevance to query  $q$ .
- ❑ Document-at-a-time scoring makes heavy use of disk seeks. With increasing query length  $|q|$ , dependent on the type of disks used, and the distribution of the index across disks, the practical run time of this approach can be poor (albeit, theoretically, exactly the same postings are processed as for term-at-a-time scoring).
- ❑ Document-at-a-time scoring has a rather small memory footprint on the order of the number of documents to return. This footprint can easily be bounded within top- $k$  retrieval by limiting the size of the results priority queue to the  $k$  entries with the currently highest scores.
- ❑ Document-at-a-time scoring presumes a global postlist ordering by document identifier or document quality.

# Query Processing II

## Disjunctive Queries

Algorithm: Term-at-a-time Scoring.

Input:  $L_1, \dots, L_m$ . The postlists of the terms  $t_1, \dots, t_m$  of query  $q$ .  
 $\mathbf{q}$ . Representation of query  $q$ , e.g., as array of  $m$  term weights.

Output: A list of documents in  $D$ , sorted in descending order of relevance to  $q$ .

*TAATScoring*( $L_1, \dots, L_m, \mathbf{q}$ )

- |   |  |
|---|--|
| 1. $R = \text{map}()$                           | 1. Initialization of result list $R$ as map.                     |
| 2. <b>FOR</b> $i \in [1, m]$ <b>DO</b>          | 2. Process postlists iteratively.                                |
| 3. $x_i = L_i.\text{head}$                      | 3. Initialization of postlist iterator for the $i$ -th postlist. |
| 4. <b>WHILE</b> $x_i \neq \text{NIL}$ <b>DO</b> | 4. For each document $d$ 's posting in the postlist:             |
| 5. $d = x_i.\text{key}$                         | 5. Get $d$ 's ID.  |
| 6. $w = x_i.\text{weight}$                      | 6. Get $t$ 's term weight for $d$ .                              |
| 7. $R[d] = R[d] + \mathbf{q}[i] \cdot w$        | 7. Update $d$ 's partial document score.                         |
| 8. $x_i = x_i.\text{next}$                      | 8. Advance the iterator.   |
| 9. <b>ENDDO</b>                                 |  |
| 10. <b>ENDDO</b>                                |  |
| 11. $\text{return}(\text{PriorityQueue}(R))$    | 11. Return the result list, ordered by document scores.          |



## Remarks:

- ❑ TAAT = Term at a time
- ❑ Term-at-a-time scoring has a comparably high main memory load, since the last “intermediate”  $|R| = |\bigcup_{i=1}^m L_i|$  before an actual ordering is performed. Otherwise, postlists are read consecutively, which suits rotating hard disks. Massive parallelization is possible.
- ❑ The order in which terms are processed (Line 2) affects how quick the intermediate scores in  $R$  approach the final document scores.
- ❑ The relevance function  $\rho$  must be additive (Line 7), or otherwise incrementally computable.
- ❑ Term-at-a-time scoring makes no a priori assumptions about postlist ordering; in case of conjunctive interpretation some ordering by document identifier is still very helpful since then skip lists can be exploited. However, to speed up retrieval and allow for (unsafe) early termination, ordering by term weight is required.

# Query Processing II

## Top-k Retrieval

Search engine users are often interested only in the top ranked  $k$  documents. Lower-ranked documents will likely never be viewed.

Query processing optimization approaches:

- ❑ **Term weight threshold**

TAAT-scoring: skip query terms whose inverse document frequency is lower than that of other query terms. Exception: stop word-heavy queries (e.g., to be or not to be).

- ❑ **Relevance score threshold**

DAAT-scoring: once  $> k$  documents have been found, determine co-occurring query terms in the top  $k$  ones; skip remaining documents not containing co-occurring query terms.

- ❑ **Early termination**

Postlists ordered by term weight: stop postlist traversal early, disregarding the rest of the postlist that cannot contribute enough to a document's relevance score.

- ❑ **Tiered indexes**

Divide documents into index tiers by quality or term frequency. If an insufficient amount of documents is found in the top tier, resort to the next one.

# Query Processing II

## Index Distribution

The larger the size of the document collection  $D$  to be indexed, the more query processing time can be improved by scaling up and scaling out.

### Term distribution

- ❑ Distribution of postlists across local disks.
- ❑ Speeds up processing on spinning hard drives.

### Document distribution (also: sharding)

- ❑ Random division of the document collection into subsets (so-called shards) and indexing of each shard on a different server for parallel query processing.
- ❑ Benefit: Smaller indexes return (more) results faster due to shorter postlists.
- ❑ Overhead: Query broker to dispatch queries and fuse each server's results.

### Tiered indexes

- ❑ Sharding of the document collection into tiers (e.g., by document importance)
- ❑ For instance: Tier 1 shards are kept in RAM, Tier 2 shards are kept in flash memory, and Tier 3 shards on spinning hard disks.

# Query Processing II

## Caching

Queries obey Zipf's law: roughly half the queries a day are unique on that day. Moreover, about 15% of the queries per day have never occurred before [\[Gomes 2017\]](#).

Consequently, the majority of queries have been seen before, enabling the use of caching to speed up query processing.

Caching can be applied at various points:

- ❑ Result caching
- ❑ Caching of postlist intersections
- ❑ Postlist caching

Individual cache refresh strategies must be employed to avoid stale data. Cache hierarchies of hardware and operating system should be exploited.

# Chapter IR:II

## II. Indexing

- ☐ Indexing Basics
- ☐ Inverted Index
- ☐ Query Processing I
- ☐ Query Processing II
- ☐ **Index Construction**
- ☐ Index Compression
- ☐ **Size Estimation**

# Index Construction

## In-Memory Index Construction

Algorithm: Index Construction.

Input:  $D = \{d_1, \dots, d_n\}$ . Set of documents  $d_i = (t_1, \dots, t_m)$  as lists of terms.

Output: Inverted index of  $D$ ; postlist of term  $t_j$  contains postings  $\boxed{i, \textcolor{violet}{tf}(t_j, d_i)}$ .

### *InMemoryIndex*( $D$ )

```
1.  $I = \text{map}()$ 
2. FOR  $i \in [1, n]$  DO
3.    $d_i = D[i]$ ;  $T = \text{set}()$ ;  $\textcolor{violet}{TF} = \text{map}()$ 
4.   FOR  $t \in d_i$  DO
5.      $\text{Insert}(T, t)$ 
6.      $\textcolor{violet}{TF}[t] = \textcolor{violet}{TF}[t] + 1$ 
7.   ENDDO
8.   FOR  $t \in T$  DO
9.     IF  $t \notin I$  THEN  $I[t] = \text{list}()$  ENDIF
10.     $L = I[t]$ 
11.     $\text{posting} = \text{record}(i, \textcolor{violet}{TF}[t])$ 
12.     $\textcolor{green}{InsertEnd}(L, \text{posting})$ 
13.   ENDDO
14. ENDDO
15.  $\text{return}(I)$ 
```

1. Initialization of an empty map as index  $I$ .
2. For each document  $d$  in  $D$ :
3. Collect  $d$ 's terminology in  $T$ .
4.  $\textcolor{violet}{Accumulate}$   $d$ 's term frequencies.
5. For each term  $t$  in  $d$ 's terminology  $T$ :
6.  $\textcolor{green}{Insert}$  new posting  $\boxed{i, \textcolor{violet}{tf}(t_j, d_i)}$  for  $d$  in  $I$ .
7. Return index of  $D$ .

# Index Construction

## Index Merging

If the document collection  $D$  does not fit into main memory, indexing is done iteratively, sharding the document collection and merging the shard indexes similar to an external merge sort:

1. The *InMemoryIndex* procedure runs until main memory is full.
2. The postlists are written to disk in alphabetical order of terms.
3. Steps 1 and 2 are repeated until  $D$  is processed.
4. All  $k$  postlist files created are read concurrently, performing a  **$k$ -way merge**.

# Index Construction

## Index Merging

$T$	Postings				
$\vdots$					
$t_i$	<div>4, 9<div><div></div><div></div><div></div></div></div> <div>19, 7<div><div></div><div></div><div></div></div></div> <div>23, 5<div><div></div><div></div><div></div></div></div> <div>28, 6<div><div></div><div></div><div></div></div></div> <div>50, 6<div><div></div><div></div><div></div></div></div> ...				
$t_j$	<div>1, 1<div><div></div><div></div><div></div></div></div> <div>3, 5<div><div></div><div></div><div></div></div></div> <div>51, 5<div><div></div><div></div><div></div></div></div> <div>60, 5<div><div></div><div></div><div></div></div></div> <div>71, 3<div><div></div><div></div><div></div></div></div> ...				
$\vdots$					

$T$	Postings					
$\vdots$						
$t_i$	<div>2, 4<div><div></div><div></div></div></div>	<div>8, 2<div><div></div><div></div></div></div>	<div>16, 1<div><div></div><div></div></div></div>	<div>41, 8<div><div></div><div></div></div></div>	<div>77, 8<div><div></div><div></div></div></div>	...
$t_j$	<div>2, 3<div><div></div><div></div></div></div>	<div>5, 2<div><div></div><div></div></div></div>	<div>8, 17<div><div></div><div></div></div></div>	<div>41, 6<div><div></div><div></div></div></div>	<div>77, 2<div><div></div><div></div></div></div>	...
$\vdots$						



# Index Construction

## Index Merging

$T$	Postings								
$\vdots$	$x_{i_1}$								
$t_i$	<table><tr><td>4, 9</td><td></td><td>19, 7</td><td>23, 5</td><td>28, 6</td><td></td><td>50, 6</td><td>...</td></tr></table>	4, 9		19, 7	23, 5	28, 6		50, 6	...
4, 9		19, 7	23, 5	28, 6		50, 6	...		
$t_j$	<table><tr><td>1, 1</td><td></td><td>3, 5</td><td>51, 5</td><td>60, 5</td><td></td><td>71, 3</td><td>...</td></tr></table>	1, 1		3, 5	51, 5	60, 5		71, 3	...
1, 1		3, 5	51, 5	60, 5		71, 3	...		
$\vdots$									

$T$	Postings								
$\vdots$	$x_{i_2}$								
$t_i$	<table><tr><td>2, 4</td><td></td><td>8, 2</td><td>16, 1</td><td>41, 8</td><td></td><td>77, 8</td><td>...</td></tr></table>	2, 4		8, 2	16, 1	41, 8		77, 8	...
2, 4		8, 2	16, 1	41, 8		77, 8	...		
$t_j$	<table><tr><td>2, 3</td><td></td><td>5, 2</td><td>8, 17</td><td>41, 6</td><td></td><td>77, 2</td><td>...</td></tr></table>	2, 3		5, 2	8, 17	41, 6		77, 2	...
2, 3		5, 2	8, 17	41, 6		77, 2	...		
$\vdots$									

$T$	Postings
$\vdots$	
$t_i$	

# Index Construction

## Index Merging

$T$	Postings								
$\vdots$									
	$x_{i_1}$								
$t_i$	<table><tr><td>4, 9</td><td></td><td>19, 7</td><td>23, 5</td><td>28, 6</td><td></td><td>50, 6</td><td>...</td></tr></table>	4, 9		19, 7	23, 5	28, 6		50, 6	...
4, 9		19, 7	23, 5	28, 6		50, 6	...		
$t_j$	<table><tr><td>1, 1</td><td></td><td>3, 5</td><td>51, 5</td><td>60, 5</td><td></td><td>71, 3</td><td>...</td></tr></table>	1, 1		3, 5	51, 5	60, 5		71, 3	...
1, 1		3, 5	51, 5	60, 5		71, 3	...		
$\vdots$									

$T$	Postings								
$\vdots$									
	$x_{i_2}$								
$t_i$	<table><tr><td>2, 4</td><td></td><td>8, 2</td><td>16, 1</td><td>41, 8</td><td></td><td>77, 8</td><td>...</td></tr></table>	2, 4		8, 2	16, 1	41, 8		77, 8	...
2, 4		8, 2	16, 1	41, 8		77, 8	...		
$t_j$	<table><tr><td>2, 3</td><td></td><td>5, 2</td><td>8, 17</td><td>41, 6</td><td></td><td>77, 2</td><td>...</td></tr></table>	2, 3		5, 2	8, 17	41, 6		77, 2	...
2, 3		5, 2	8, 17	41, 6		77, 2	...		
$\vdots$									

$T$	Postings		
$\vdots$			
$t_i$	<table><tr><td>2, 4</td><td></td></tr></table>	2, 4	
2, 4			

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	<div>4, 9<div><div></div><div></div><div></div></div></div>	<div>19, 7<div><div></div><div></div><div></div></div></div>	<div>23, 5<div><div></div><div></div><div></div></div></div>	<div>28, 6<div><div></div><div></div><div></div></div></div>	<div>50, 6<div><div></div><div></div><div></div></div></div>	$\dots$
$t_j$	<div>1, 1<div><div></div><div></div><div></div></div></div>	<div>3, 5<div><div></div><div></div><div></div></div></div>	<div>51, 5<div><div></div><div></div><div></div></div></div>	<div>60, 5<div><div></div><div></div><div></div></div></div>	<div>71, 3<div><div></div><div></div><div></div></div></div>	$\dots$
$\vdots$						

$T$	Postings				
$\vdots$	$x_{i_2}$				
$t_i$	<div><div>2, 4</div><div></div><div></div></div>	<div><div>8, 2</div></div>	<div><div>16, 1</div></div>	<div><div>41, 8</div><div></div></div>	<div><div>77, 8</div></div> ...
$t_j$	<div><div>2, 3</div><div></div><div></div></div>	<div><div>5, 2</div></div>	<div><div>8, 17</div></div>	<div><div>41, 6</div><div></div></div>	<div><div>77, 2</div></div> ...
$\vdots$					

$T$	Postings	
$\vdots$		
$t_i$	2, 4	4, 9

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	<div>4, 9<div><div></div><div></div><div></div></div></div>	<div>19, 7<div><div></div><div></div><div></div></div></div>	<div>23, 5<div><div></div><div></div><div></div></div></div>	<div>28, 6<div><div></div><div></div><div></div></div></div>	<div>50, 6<div><div></div><div></div><div></div></div></div>	$\dots$
$t_j$	<div>1, 1<div><div></div><div></div><div></div></div></div>	<div>3, 5<div><div></div><div></div><div></div></div></div>	<div>51, 5<div><div></div><div></div><div></div></div></div>	<div>60, 5<div><div></div><div></div><div></div></div></div>	<div>71, 3<div><div></div><div></div><div></div></div></div>	$\dots$
$\vdots$						

$T$	Postings				
$\vdots$	$x_{i_2}$				
$t_i$	<div>2, 4<div><div></div><div></div><div></div></div></div> <div>8, 2</div> <div>16, 1</div> <div>41, 8<div><div></div><div></div><div></div></div></div> <div>77, 8</div> ...				
$t_j$	<div>2, 3<div><div></div><div></div><div></div></div></div> <div>5, 2</div> <div>8, 17</div> <div>41, 6<div><div></div><div></div><div></div></div></div> <div>77, 2</div> ...				
$\vdots$					

$T$	Postings		
$\vdots$			
$t_i$	<div>2, 4</div>	<div></div>	<div>4, 9</div>

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	<div>4, 9<div><div></div><div></div><div></div></div></div>	<div>19, 7<div><div></div><div></div><div></div></div></div>	<div>23, 5<div><div></div><div></div><div></div></div></div>	<div>28, 6<div><div></div><div></div><div></div></div></div>	<div>50, 6<div><div></div><div></div><div></div></div></div>	$\dots$
$t_j$	<div>1, 1<div><div></div><div></div><div></div></div></div>	<div>3, 5<div><div></div><div></div><div></div></div></div>	<div>51, 5<div><div></div><div></div><div></div></div></div>	<div>60, 5<div><div></div><div></div><div></div></div></div>	<div>71, 3<div><div></div><div></div><div></div></div></div>	$\dots$
$\vdots$						

$T$	Postings					
$\vdots$	$x_{i_2}$					
$t_i$	<div>2, 4<div><div></div><div></div></div></div>	<div>8, 2<div><div></div><div></div></div></div>	<div>16, 1<div><div></div><div></div></div></div>	<div>41, 8<div><div></div><div></div></div></div>	<div>77, 8<div><div></div><div></div></div></div>	$\dots$
$t_j$	<div>2, 3<div><div></div><div></div></div></div>	<div>5, 2<div><div></div><div></div></div></div>	<div>8, 17<div><div></div><div></div></div></div>	<div>41, 6<div><div></div><div></div></div></div>	<div>77, 2<div><div></div><div></div></div></div>	$\dots$
$\vdots$						

$T$	Postings			
$\vdots$				
$t_i$	<div>2, 4<div><div></div><div></div></div></div>	<div>4, 9<div><div></div><div></div></div></div>	<div>8, 2<div><div></div><div></div></div></div>	<div>16, 1<div><div></div><div></div></div></div>

# Index Construction

## Index Merging

$T$	Postings				
$\vdots$	$x_{i_1}$				
$t_i$	<div>4, 9<div><div></div><div></div></div></div> <div>19, 7</div> <div>23, 5</div> <div>28, 6<div><div></div><div></div></div></div> <div>50, 6</div> <div>...</div>				
$t_j$	<div>1, 1<div><div></div><div></div></div></div> <div>3, 5</div> <div>51, 5</div> <div>60, 5<div><div></div><div></div></div></div> <div>71, 3</div> <div>...</div>				
$\vdots$					

$T$	Postings				
$\vdots$	$x_{i_2}$				
$t_i$	<div><div>2, 4</div><div></div><div></div></div>	<div><div>8, 2</div><div></div><div></div></div>	<div><div>16, 1</div><div></div><div></div></div>	<div><div>41, 8</div><div></div><div></div></div>	<div><div>77, 8</div><div></div><div></div></div> ...
$t_j$	<div><div>2, 3</div><div></div><div></div></div>	<div><div>5, 2</div><div></div><div></div></div>	<div><div>8, 17</div><div></div><div></div></div>	<div><div>41, 6</div><div></div><div></div></div>	<div><div>77, 2</div><div></div><div></div></div> ...
$\vdots$					

$T$	Postings					
$\vdots$						
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	4, 9		19, 7	23, 5	28, 6	50, 6 ...
$t_j$	1, 1		3, 5	51, 5	60, 5	71, 3 ...
$\vdots$						

$T$	Postings					
$\vdots$	$x_{i_2}$					
$t_i$	2, 4		8, 2	16, 1	41, 8	77, 8 ...
$t_j$	2, 3		5, 2	8, 17	41, 6	77, 2 ...
$\vdots$						

$T$	Postings					
$\vdots$						
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7 23, 5

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	4, 9		19, 7	23, 5	28, 6	50, 6 ...
$t_j$	1, 1		3, 5	51, 5	60, 5	71, 3 ...
$\vdots$						

$T$	Postings					
$\vdots$	$x_{i_2}$					
$t_i$	2, 4		8, 2	16, 1	41, 8	77, 8 ...
$t_j$	2, 3		5, 2	8, 17	41, 6	77, 2 ...
$\vdots$						

$T$	Postings							
$\vdots$								
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7	23, 5	28, 6



# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	4, 9		19, 7	23, 5	28, 6	50, 6 ...
$t_j$	1, 1		3, 5	51, 5	60, 5	71, 3 ...
$\vdots$						

$T$	Postings					
$\vdots$	$x_{i_2}$					
$t_i$	2, 4		8, 2	16, 1	41, 8	77, 8 ...
$t_j$	2, 3		5, 2	8, 17	41, 6	77, 2 ...
$\vdots$						

$T$	Postings							
$\vdots$								
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7	23, 5	28, 6 41, 8

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	4, 9		19, 7	23, 5	28, 6	50, 6 ...
$t_j$	1, 1		3, 5	51, 5	60, 5	71, 3 ...
$\vdots$						

$T$	Postings					
$\vdots$	$x_{i_2}$					
$t_i$	2, 4		8, 2	16, 1	41, 8	77, 8 ...
$t_j$	2, 3		5, 2	8, 17	41, 6	77, 2 ...
$\vdots$						

$T$	Postings									
$\vdots$										
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7	23, 5	28, 6	41, 8	50, 6

# Index Construction

## Index Merging

$T$	Postings					
$\vdots$	$x_{i_1}$					
$t_i$	4, 9		19, 7	23, 5	28, 6	50, 6 ...
$t_j$	1, 1		3, 5	51, 5	60, 5	71, 3 ...
$\vdots$						

$T$	Postings					
$\vdots$	$x_{i_2}$					
$t_i$	2, 4		8, 2	16, 1	41, 8	77, 8 ...
$t_j$	2, 3		5, 2	8, 17	41, 6	77, 2 ...
$\vdots$						

$T$	Postings											
$\vdots$												
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7	23, 5	28, 6	41, 8	50, 6	77, 8	...

# Index Construction

## Index Merging

<i>T</i>	Postings					
⋮						
$t_i$	4, 9		19, 7	23, 5	28, 6	50, 6 ...
$t_j$	1, 1		3, 5	51, 5	60, 5	71, 3 ...
⋮						

<i>T</i>	Postings					
⋮						
$t_i$	2, 4		8, 2	16, 1	41, 8	77, 8 ...
$t_j$	2, 3		5, 2	8, 17	41, 6	77, 2 ...
⋮						

<i>T</i>	Postings											
⋮												
$t_i$	2, 4		4, 9	8, 2	16, 1	19, 7	23, 5	28, 6	41, 8	50, 6	77, 8	...
$t_j$	1, 1		2, 3	3, 5	5, 2	8, 17	41, 6	51, 5	60, 5	71, 3	77, 2	...
⋮												

## Remarks:

- ❑ Alphabetical ordering of intermediary postlist files ensures that the index can be read sequentially, albeit concurrently, during merging. Compare to [document-at-a-time scoring](#).
- ❑ If a term appears in only one of the indexes, its postlist is directly added to the merged index.
- ❑ Postings with skip pointers can be pre-determined before merging a postlist so that appropriate space can be allocated immediately, but the actual skip pointers need to be recomputed after the postlist is merged.
- ❑ The number  $k$  of intermediary postlist files that can be read concurrently without causing too much seeking overhead depends on the underlying hardware (e.g.,  $k$  is smaller for spinning hard disks than for solid state disks). In case  $k$  is too large, the intermediary postlist files are merged in multiple passes,  $k' < k$  at a time, until all are merged.

# Index Construction

## Distributed Indexing

If neither the document collection  $D$ , nor its index can be stored on a single machine, indexing must be performed distributed across a computer cluster.

Many cluster computing frameworks exist nowadays; the [NoSQL movement](#), and ultimately the Big Data hype, was kicked off by Google's MapReduce [[Dean 2004](#)].

# Index Construction

## Distributed Indexing

If neither the document collection  $D$ , nor its index can be stored on a single machine, indexing must be performed distributed across a computer cluster.

Many cluster computing frameworks exist nowadays; the [NoSQL movement](#), and ultimately the Big Data hype, was kicked off by Google's MapReduce [\[Dean 2004\]](#).

From a developer perspective, data processing with MapReduce boils down to implementing two procedures:

- ❑ Map: Given a key-value pair as input, it outputs a list of key-values pairs.
- ❑ Reduce: Given a key and the list of values output by map under that key, it outputs a key-value pair.

# Index Construction

## Distributed Indexing

If neither the document collection  $D$ , nor its index can be stored on a single machine, indexing must be performed distributed across a computer cluster.

Many cluster computing frameworks exist nowadays; the [NoSQL movement](#), and ultimately the Big Data hype, was kicked off by Google's MapReduce [\[Dean 2004\]](#).

From a developer perspective, data processing with MapReduce boils down to implementing two procedures:

- ❑ Map: Given a key-value pair as input, it outputs a list of key-values pairs.
- ❑ Reduce: Given a key and the list of values output by map under that key, it outputs a key-value pair.

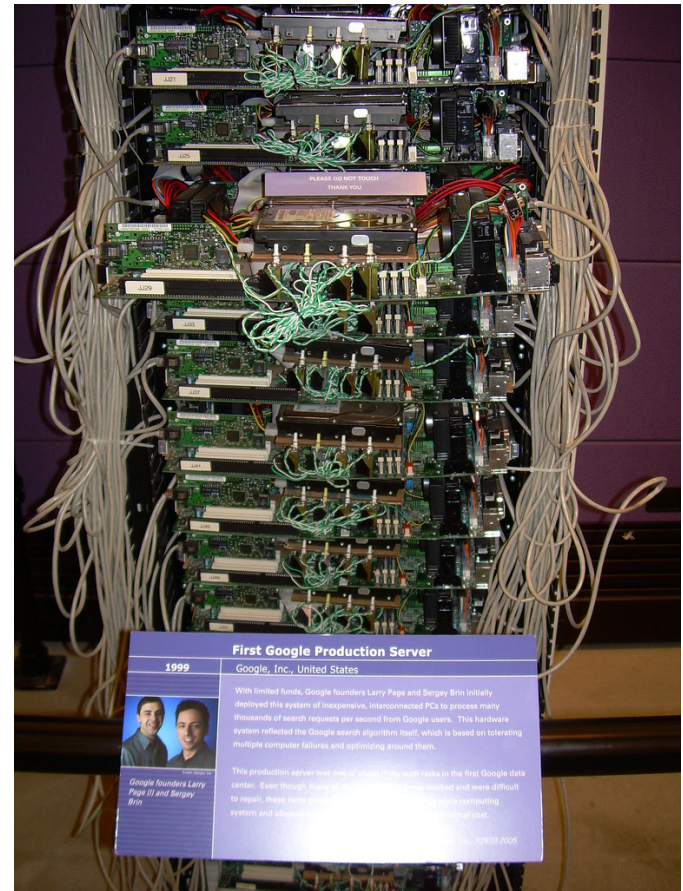
Example:

- ❑ IndexingMapper: Given a pair  $(i, d_i)$ , where  $i$  is the document identifier of document  $d_i$  as input, output a pair  $(t, i)$  for every unique  $t \in d_i$ .
- ❑ **DF**Reducer: Given  $(t, [\dots, i, \dots])$  as input, output  $(t, |[\dots, i, \dots]|) = (t, \text{df}(t, D))$ .



## Remarks:

- ❑ Computer clusters are often built from inexpensive commodity hardware. In the early days, desktop computers were used as [Beowulf clusters](#), or dismantled and stacked. Google 1997 and 1999:



- ❑ The key contributions of the MapReduce framework are not the actual map and reduce functions, but the scalability and fault-tolerance achieved for a variety of applications by optimizing the execution engine [\[Wikipedia\]](#).
- ❑ This framework is best-suited for problems that are [embarrassingly parallel](#).
- ❑ The most widespread open source implementation is found in [Apache Hadoop](#).

# Index Construction

## Distributed Indexing

Presuming the document collection is stored in a distributed document storage across the cluster, the execution of a MapReduce job divides into three basic phases:

- ❑ **Map phase**

The map function is called in parallel on all cluster nodes and fed chunks of the data. Its output is recorded locally on each cluster node.

- ❑ **Shuffle phase**

The output is transferred to a random cluster node chosen using a hash function, so that the same key is always transferred to the same cluster node. Once all data belonging to a key are on the same node, the values are sorted.

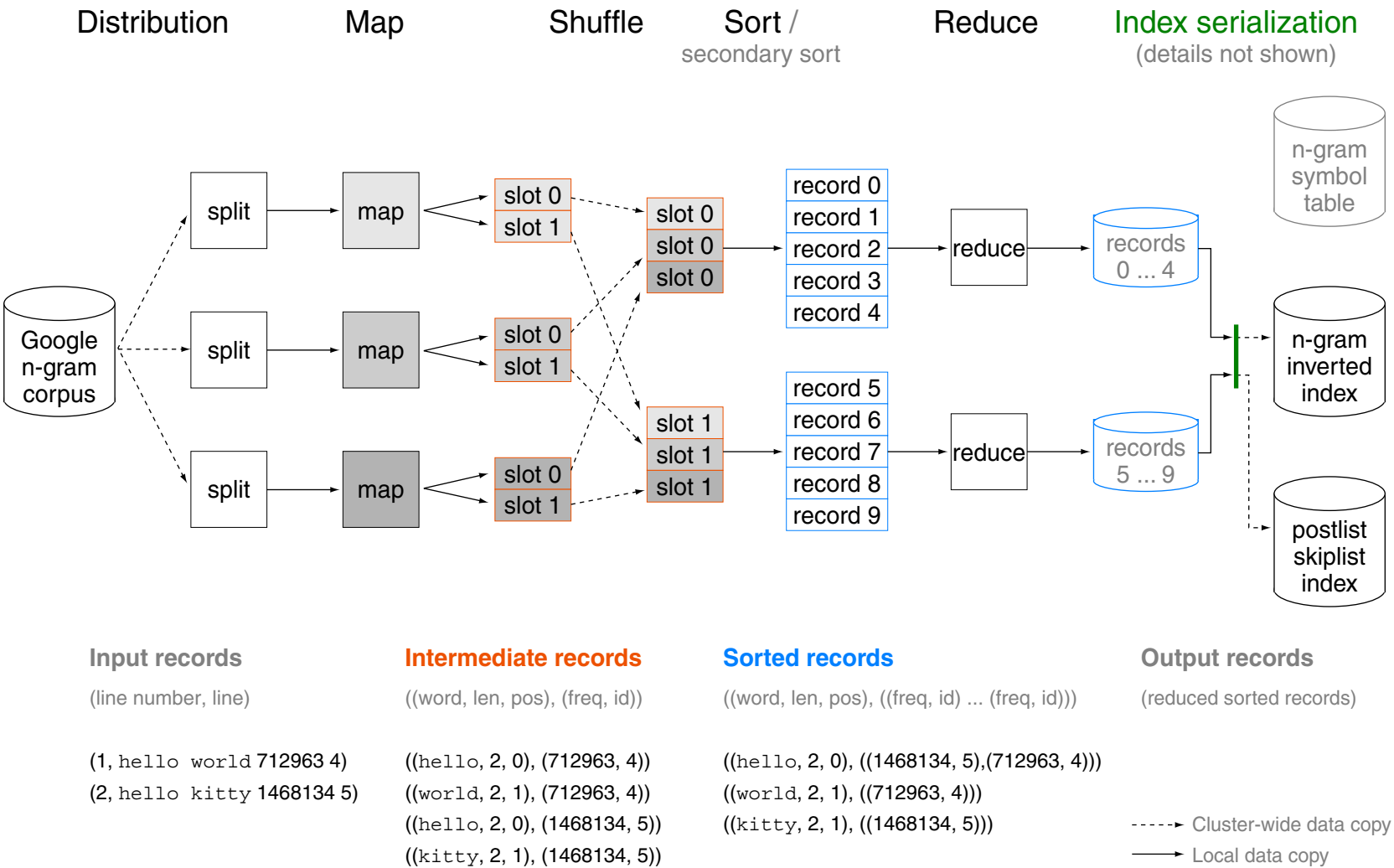
- ❑ **Reduce phase**

The reduce function is called in parallel on all cluster nodes and fed the sorted lists, recording their output.

The map and reduce functions are idempotent: they are reexecuted in case of failures. To make optimal use of available resources, the framework may execute the same task more than once on different machines, retaining the first output that emerges (so-called speculative execution).

# Index Construction

## Distributed Indexing: Example Netspeak [\[www.netspeak.org\]](http://www.netspeak.org)



# Index Construction

## Index Updates

Document collections grow and change. Therefore, the index must be updated. The following strategies are applied:

- ❑ **Index merging**

When new documents arrive in large numbers at a time, they are indexed and then the existing index is merged with the new one.

- ❑ **Result merging**

When new documents arrive in small numbers at a time, a separate, small index is maintained and updated. Queries are processed against both the existing index and the small one containing the new arrivals, fusing the results.

- ❑ **Deletions list**

Deletions are recorded in a deletions list, and deleted documents are removed from search results before results are shown.

Modifications are done by inserting a new document, and deleting the previous version.

# Chapter IR:II

## II. Indexing

- ❑ Indexing Basics
- ❑ Inverted Index
- ❑ Query Processing I
- ❑ Query Processing II
- ❑ Index Construction
- ❑ Index Compression
- ❑ Size Estimation

# Size Estimation

## Query Result Set Size

tropical fish aquarium

Search

Web results

Page 1 of 3,880,000 results

The total number of results is estimated, since web search engines typically do not explore the entire indexed document collection to compute the first page of results returned, but only a subset.

### Approaches:

- ❑ Joint probability estimation
- ❑ Conditional probability estimation
- ❑ Initial result set-based estimation

## Remarks:

- Example data from the GOV2 collection (collection size  $|D|$  is 25,205,179):

Query	Document frequency
aquarium	26,480
breeding	81,885
fish	1,131,855
lincoln	771,326
tropical	120,990
aquarium breeding	1,848
fish aquarium	9,722
fish breeding	36,427
tropical aquarium	1,921
tropical breeding	5,510
tropical fish	18,472
tropical fish aquarium	1,529
tropical fish breeding	3,629

# Size Estimation

## Query Result Set Size: Joint Probability

Let  $P_{df}(t)$  denote the probability of  $t$  occurring at least once in a document:

$$P_{df}(t) = \frac{df(t)}{|D|},$$

where  $D$  denotes the document collection of size  $|D|$  and  $df(t)$  the number of documents in  $D$  containing  $t$ , called document frequency.



# Size Estimation

## Query Result Set Size: Joint Probability

Let  $P_{df}(t)$  denote the probability of  $t$  occurring at least once in a document:

$$P_{df}(t) = \frac{df(t)}{|D|},$$

where  $D$  denotes the document collection of size  $|D|$  and  $df(t)$  the number of documents in  $D$  containing  $t$ , called document frequency.

The result set size  $df(q)$  of a query  $q$  of length  $|q|$  terms can be estimated with

$$df(q) = |D| \cdot \prod_{i=1}^{|q|} P_{df}(t_i) = \frac{\prod_{i=1}^{|q|} df(t_i)}{|D|^{|q|-1}},$$

where  $t_i$  denotes the  $i$ -th term in  $q$ . This estimation presumes term independence.

# Size Estimation

## Query Result Set Size: Joint Probability

Let  $P_{df}(t)$  denote the probability of  $t$  occurring at least once in a document:

$$P_{df}(t) = \frac{df(t)}{|D|},$$

where  $D$  denotes the document collection of size  $|D|$  and  $df(t)$  the number of documents in  $D$  containing  $t$ , called document frequency.

The result set size  $df(q)$  of a query  $q$  of length  $|q|$  terms can be estimated with

$$df(q) = |D| \cdot \prod_{i=1}^{|q|} P_{df}(t_i) = \frac{\prod_{i=1}^{|q|} df(t_i)}{|D|^{|q|-1}},$$

where  $t_i$  denotes the  $i$ -th term in  $q$ . This estimation presumes term independence.

Examples:

❑  $df(\text{tropical fish aquarium}) = 5.71$

actual: 1,529 documents

❑  $df(\text{tropical fish breeding}) = 17.65$

actual: 3,629 documents

# Size Estimation

## Query Result Set Size: Conditional Probability

By exploiting term co-occurrence information, we can obtain better estimates with

$$P_{df}(q) = P_{df}(t_1, t_2, t_3) = P_{df}(t_1, t_2) \cdot P_{df}(t_3 \mid t_1, t_2),$$

where  $P_{df}(t_3 \mid t_1, t_2) \approx P_{df}(t_3 \mid t_x) = \max\{P_{df}(t_3 \mid t_1), P_{df}(t_3 \mid t_2)\}$  and  $|q| = 3$ . Recall that  $P(A \mid B) = P(A, B)/P(B)$ . Hence

$$df(q) = |D| \cdot P_{df}(q) = \frac{df(t_1, t_2) \cdot df(t_x, t_3)}{df(t_x)}.$$

# Size Estimation

## Query Result Set Size: Conditional Probability

By exploiting term co-occurrence information, we can obtain better estimates with

$$P_{df}(q) = P_{df}(t_1, t_2, t_3) = P_{df}(t_1, t_2) \cdot P_{df}(t_3 \mid t_1, t_2),$$

where  $P_{df}(t_3 \mid t_1, t_2) \approx P_{df}(t_3 \mid t_x) = \max\{P_{df}(t_3 \mid t_1), P_{df}(t_3 \mid t_2)\}$  and  $|q| = 3$ . Recall that  $P(A \mid B) = P(A, B)/P(B)$ . Hence

$$df(q) = |D| \cdot P_{df}(q) = \frac{df(t_1, t_2) \cdot df(t_x, t_3)}{df(t_x)}.$$

- ❑ Queries of length  $|q| = 2$  need not be estimated, anymore.
- ❑ Queries of length  $|q| = 3$  are typically underestimated.
- ❑ Queries of length  $|q| > 3$  still require estimations based on term independence, or storing higher-order co-occurrence information.

Examples:

- ❑  $df(\text{tropical fish aquarium}) = 293$  actual: 1,529 documents
- ❑  $df(\text{tropical fish breeding}) = 841$  actual: 3,629 documents

# Size Estimation

## Query Result Set Size: Initial Result Set-based Estimation

Let  $D' \subset D$  denote the documents **initially** scored for a query  $q$  (e.g., in tiered index). Then the size of the total result set in  $D$  can be estimated with

$$df(q) = |D_t| \cdot \frac{|D'_q|}{|D'|} = df(t) \cdot \frac{|\{d \mid d \in D' \wedge q \in d\}|}{|D'|},$$

where  $t$  is the query term with the smallest subset  $D_t \subset D$  of documents that contain  $t$ , and  $D'_q \subset D'$  is the subset of the initially scored documents  $D'$  that contain all terms of  $q$ .

# Size Estimation

## Query Result Set Size: Initial Result Set-based Estimation

Let  $D' \subset D$  denote the documents **initially** scored for a query  $q$  (e.g., in tiered index). Then the size of the total result set in  $D$  can be estimated with

$$df(q) = |D_t| \cdot \frac{|D'_q|}{|D'|} = df(t) \cdot \frac{|\{d \mid d \in D' \wedge q \in d\}|}{|D'|},$$

where  $t$  is the query term with the smallest subset  $D_t \subset D$  of documents that contain  $t$ , and  $D'_q \subset D'$  is the subset of the initially scored documents  $D'$  that contain all terms of  $q$ .

This estimation presumes relevant documents are uniformly distributed across all documents in  $D_t$ . **Why does it usually overestimate the result set size?**

Examples:

- ❑ With  $D_{\text{aquarium}} = 26,480$ , let  $|D'| = 3,000$ , and  $|D'_q| = 258$ :  
❑  $df(\text{tropical fish aquarium}) = 2,277$  actual: 1,529 documents
- ❑ With  $D_{\text{breeding}} = 81,885$ , let  $|D'| = 3,000$ , and  $|D'_q| = 150$ :  
❑  $df(\text{tropical fish breeding}) = 4,094$  actual: 3,629 documents

# Size Estimation

## Query Result Set Size: Initial Result Set-based Estimation

Let  $D' \subset D$  denote the documents **initially** scored for a query  $q$  (e.g., in tiered index). Then the size of the total result set in  $D$  can be estimated with

$$df(q) = |D_t| \cdot \frac{|D'_q|}{|D'|} = df(t) \cdot \frac{|\{d \mid d \in D' \wedge q \in d\}|}{|D'|},$$

where  $t$  is the query term with the smallest subset  $D_t \subset D$  of documents that contain  $t$ , and  $D'_q \subset D'$  is the subset of the initially scored documents  $D'$  that contain all terms of  $q$ .

This estimation presumes relevant documents are uniformly distributed across all documents in  $D_t$ . Overestimations result from  $D'$  containing the most “important” documents indexed. As  $|D'|$  **approaches**  $|D_w|$ , estimations approach the true value.

Examples:

□ With  $D_{\text{aquarium}} = 26,480$ , let  $|D'| = 6,000$ , and  $|D'_q| = 402$ :

□  $df(\text{tropical fish aquarium}) = 1,774$  actual: 1,529 documents

□ With  $D_{\text{breeding}} = 81,885$ , let  $|D'| = 6,000$ , and  $|D'_q| = 276$ :

□  $df(\text{tropical fish breeding}) = 3,767$  actual: 3,629 documents

# Size Estimation

## Indexed Collection Size: Joint Probability-based

Most search engines are black boxes to outsiders, and many do not share the size of the document collection they index.



# Size Estimation

## Indexed Collection Size: Joint Probability-based

Most search engines are black boxes to outsiders, and many do not share the size of the document collection they index.

Given a web search engine, the size  $|D|$  of the document collection  $D$  indexed can be estimated using two **independently** occurring terms  $t_1$  and  $t_2$ :

$$P_{df}(t_1, t_2) = P_{df}(t_1) \cdot P_{df}(t_2) \quad \leadsto \quad |D| = \frac{df(t_1) \cdot df(t_2)}{df(t_1, t_2)}.$$

Averaging over many term pairs improves the estimate.

Example for GOV2:

- $df(\text{tropical}) = 120,990$ ,  
 $df(\text{lincoln}) = 771,326$ , and  
 $df(\text{tropical}, \text{lincoln}) = 3018$ .

- Then  $|D| = 30,922,045$

actual: 25,205,179 documents

# Size Estimation

## Indexed Collection Size: Proportionality [\[van den Bosch 2016\]](#) [\[worldwidewebsize.com\]](http://worldwidewebsize.com)

Given a web search engine, the size  $|D|$  of the document collection  $D$  indexed can be estimated when presuming proportionality to a different reference collection  $D'$ :

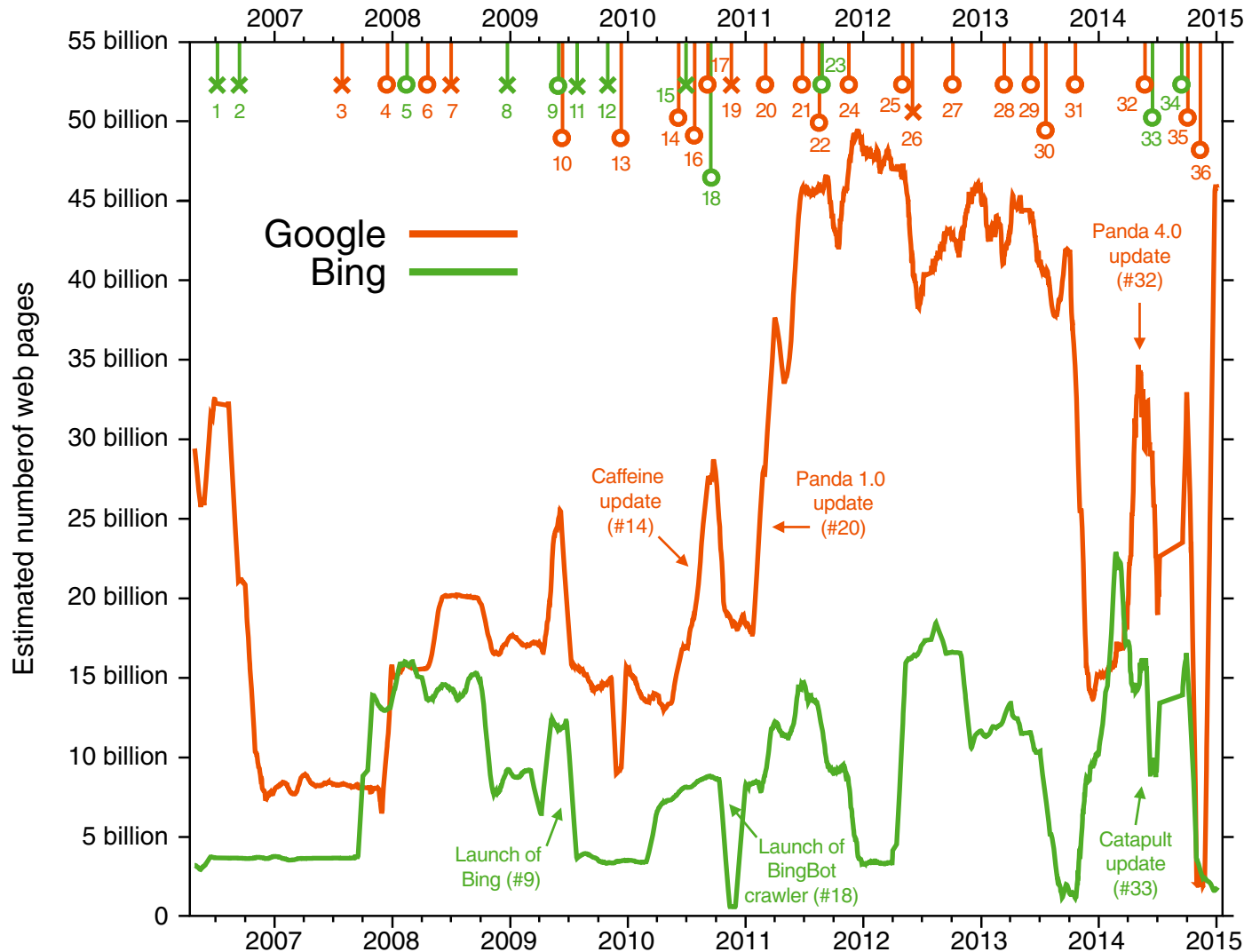
$$P_{df}^{(D)}(t) = P_{df}^{(D')}(t) \quad \leadsto \quad |D| = \frac{df_D(t) \cdot |D'|}{df_{D'}(t)},$$

where  $t$  is a term occurring in both  $D$  and  $D'$ , and  $df_D$  ( $df_{D'}$ ) computes the document frequency for  $D$  ( $D'$ ).

Averaging over terms of varying frequencies improves the estimate.

# Size Estimation

Indexed Collection Size: Proportionality [van den Bosch 2016] [worldwidewebsite.com]



## Remarks:

1. [\[2006-07-04\]](#) MSN Search outage
2. [\[2006-09-11\]](#) Launch of (improvements to) Live Search
3. [\[2007-07-31\]](#) Update to supplemental results indexing
4. [\[2007-12-18\]](#) No more supplemental index; whole index is searched for every query
5. [\[2008-02-12\]](#) Crawler improvements for Live Search
6. [\[2008-04-11\]](#) Improved crawling of HTML forms
7. [\[2008-06-30\]](#) Improved Flash indexing
8. [\[2008-12-11\]](#) First experiments with MSNBot 2.0
9. [\[2009-05-28\]](#) Launch of Bing
10. [\[2009-06-18\]](#) Improved Flash indexing
11. [\[2009-07-31\]](#) Bing and Yahoo! team up on search
12. [\[2009-11-04\]](#) MSNBot 2.0
13. [\[2009-12-07\]](#) Updates to real-time search
14. [\[2010-06-08\]](#) Launch of new web indexing system Caffeine
15. [\[2010-06-28\]](#) Experiments with BingBot crawler
16. [\[2010-07-29\]](#) Improved Flash & AJAX indexing
17. [\[2010-08-31\]](#) Google indexes SVG
18. [\[2010-09-03\]](#) Launch of BingBot crawler
19. [\[2010-11-11\]](#) Improved Flash indexing
20. [\[2011-02-24\]](#) Panda Refresh (update to promote (English) high-quality sites more)

## Remarks: (continued)

21. [\[2011-06-21\]](#) Panda 2.2
22. [\[2011-08-12\]](#) Panda (rolled out to all languages)
23. [\[2011-08-15\]](#) Gradual roll-out of Tiger indexing architecture
24. [\[2011-11-03\]](#) Panda (update, affects 35% of queries)
25. [\[2012-04-24\]](#) Penguin update (targeting Web spam, impacting around 3.1% of queries)
26. [\[2012-05-26\]](#) Penguin 2 update (impacting less than 0.1% of queries)
27. [\[2012-10-05\]](#) Penguin 3 update (impacting around 0.3% of queries)
28. [\[2013-03-12\]](#) Panda update
29. [\[2013-05-22\]](#) Penguin 4 (v2.0, impacting 2.3% of queries)
30. [\[2013-07-18\]](#) Panda update
31. [\[2013-10-04\]](#) Penguin 5 (v2.1, impacting around 1% of queries)
32. [\[2014-05-21\]](#) Panda 4.0
33. [\[2014-06-18\]](#) Launch of Bing Catapult
34. [\[2014-09-09\]](#) Improved spam filtering
35. [\[2014-09-26\]](#) Panda 4.1 (3-5% of queries affected)
36. [\[2014-10-17\]](#) Penguin 6 (v3.0, impacting less than 1% English queries)