

# Information Retrieval (5LN712)

## Boolean Retrieval

Ali Basirat

Department of Linguistics and Philology  
Uppsala University

March 30, 2020

- 1 Boolean Retrieval
- 2 Inverted index
- 3 Processing Boolean Queries
- 4 Summary
- 5 Exercises
- 6 References

- Index documents with their terms
- Take a Boolean query and process it
- Return a set of relevant documents

- Construct a term document incidence matrix
- Each column is a binary vector representing a document
- Each row is a binary vector representing a term

## Example

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

- A query is what an end user conveys to the computer to communicate an information need
- A Boolean query consists of terms and Boolean (logical) operators (AND, OR, NOT, etc.) between them

---

**Algorithm 1** Finding relevant documents for a Boolean query

---

**Input:** A Boolean query  $q$  and a term-doc. incidence matrix  $M$

**Output:** A set of document ids

- 1: Select row vectors of  $M$  corresponding to the terms of  $q$
  - 2: Apply Boolean operators defined by  $q$  on the selected rows
  - 3: Return document ids whose values in the final Boolean vector is 1
-

## Example

- $q$ : Brutus AND Caesar AND NOT Calpurnia

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
• $M$ :	Antony	1	1	0	0	1
	Brutus	1	1	0	1	0
	Caesar	1	1	0	1	1
	Calpurnia	0	1	0	0	0
	Cleopatra	1	0	0	0	0
	mercy	1	0	1	1	1
	worser	1	0	1	1	0
	...					

## Example

- $q$ : Brutus AND Caesar AND NOT Calpurnia

- $M$ :

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						



## Example

- $q$ : Brutus AND Caesar AND NOT Calpurnia

- $M$ :

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

## Example

- $q$ : Brutus AND Caesar AND NOT Calpurnia

- $M$ :

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

## Example

- $q$ : Brutus AND Caesar AND NOT Calpurnia

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

- $[110100] \wedge [110111] \wedge \neg[010000] = [100100]$

## Example

- $q$ : Brutus AND Caesar AND NOT Calpurnia

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

- $[110100] \wedge [110111] \wedge \neg[010000] = [100100]$
- return** {Anthony and Cleopatra, Hamlet}

- A very large and sparse matrix
- Only ones are stored
- Terms are associated with the documents they occur in (inverted index)

- 1 Boolean Retrieval
- 2 Inverted index**
- 3 Processing Boolean Queries
- 4 Summary
- 5 Exercises
- 6 References

Each term is associated with a list, called a posting list, whose elements are the documents in which the term is occurred

Antony → Antony and Cleopatra → Julius Caesar → Macbeth

Brutus → Antony and Cleopatra → Julius Caesar → Hamlet

Caesar → Antony and Cleopatra → Julius Caesar → Hamlet → Othello → Macbeth

...

The major steps to construct an inverted index:

- ① collect the documents to be indexed
- ② tokenize the text
- ③ do linguistic pre-processing (token normalisation)
- ④ index the documents



- **Doc1:** I did enact Julius Caesar: I was killed i' the Capitol; Brutus killed me.
- **Doc2:** So let it be with Caesar. The noble Brutus hath told you Caesar was ambitious:

Ali Basirat

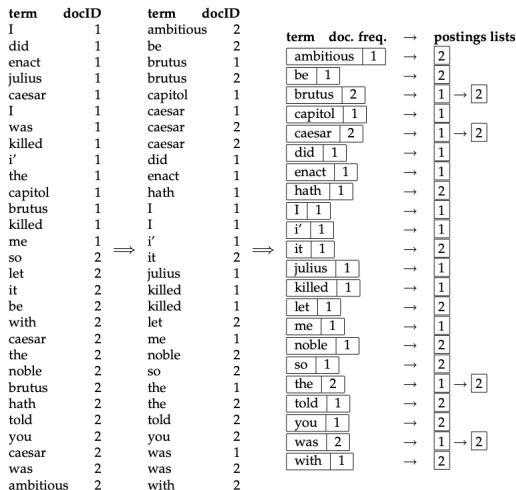


Figure: An example of the inverted index [1]

- 1 Boolean Retrieval
- 2 Inverted index
- 3 Processing Boolean Queries**
- 4 Summary
- 5 Exercises
- 6 References

The use of an inverted index instead of a term-document matrix entails a move from Boolean algebra to the set theory, which are related to each other.

## Example

- $\wedge \implies \cap$  (intersect)
- $\vee \implies \cup$  (union)
- $\neg \implies -$  (subtract)

## Definition

Conjunctive query: a query consisting of terms and logical AND ( $\wedge$ ) operator in between

## Problem

*How to process a conjunctive query using inverted indices?*

## Solution

- *locate terms in the dictionary*
- *retrieve their postings*
- *intersect the postings*

## How to intersect two postings?

```
INTERSECT( $p_1, p_2$ )
1   $answer \leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $docID(p_1) = docID(p_2)$ 
4      then  $\text{ADD}(answer, docID(p_1))$ 
5           $p_1 \leftarrow next(p_1)$ 
6           $p_2 \leftarrow next(p_2)$ 
7  else if  $docID(p_1) < docID(p_2)$ 
8      then  $p_1 \leftarrow next(p_1)$ 
9      else  $p_2 \leftarrow next(p_2)$ 
10 return  $answer$ 
```

**Figure:** The intersect algorithm. It is assumed that the postings are sorted in advance.

How to process a conjunctive query consisting of an arbitrary number of terms?

```
INTERSECT( $\langle t_1, \dots, t_n \rangle$ )  
1  $terms \leftarrow \text{SORTBYINCREASINGFREQUENCY}(\langle t_1, \dots, t_n \rangle)$   
2  $result \leftarrow \text{postings}(\text{first}(terms))$   
3  $terms \leftarrow \text{rest}(terms)$   
4 while  $terms \neq \text{NIL}$  and  $result \neq \text{NIL}$   
5 do  $result \leftarrow \text{INTERSECT}(result, \text{postings}(\text{first}(terms)))$   
6    $terms \leftarrow \text{rest}(terms)$   
7 return  $result$ 
```

Figure: The extended intersect algorithm.

## Skip-pointers in postings lists

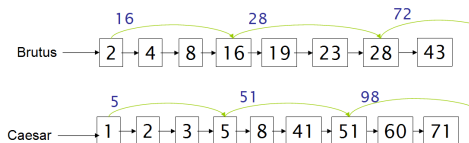


Figure: Postings lists with skip pointers



The intersection of two postings lists augmented with skip-pointers

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1   $answer \leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $docID(p_1) = docID(p_2)$ 
4      then  $\text{ADD}(answer, docID(p_1))$ 
5           $p_1 \leftarrow next(p_1)$ 
6           $p_2 \leftarrow next(p_2)$ 
7  else if  $docID(p_1) < docID(p_2)$ 
8      then if  $hasSkip(p_1)$  and  $(docID(skip(p_1)) \leq docID(p_2))$ 
9          then while  $hasSkip(p_1)$  and  $(docID(skip(p_1)) \leq docID(p_2))$ 
10             do  $p_1 \leftarrow skip(p_1)$ 
11             else  $p_1 \leftarrow next(p_1)$ 
12         else if  $hasSkip(p_2)$  and  $(docID(skip(p_2)) \leq docID(p_1))$ 
13             then while  $hasSkip(p_2)$  and  $(docID(skip(p_2)) \leq docID(p_1))$ 
14                 do  $p_2 \leftarrow skip(p_2)$ 
15                 else  $p_2 \leftarrow next(p_2)$ 
16  return  $answer$ 

```

Figure: Postings lists intersection with skip pointers.

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9  then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10     do  $p_1 \leftarrow \text{skip}(p_1)$ 
11     else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13  then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14     do  $p_2 \leftarrow \text{skip}(p_2)$ 
15     else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

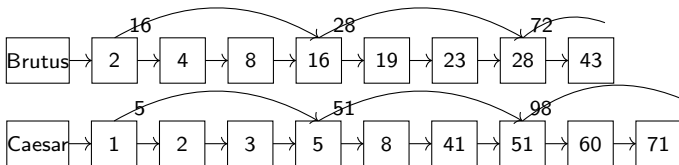


Figure: Postings lists with skip pointers

answer = { }

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

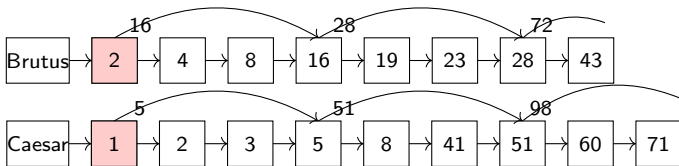


Figure: Postings lists with skip pointers

answer = { }

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

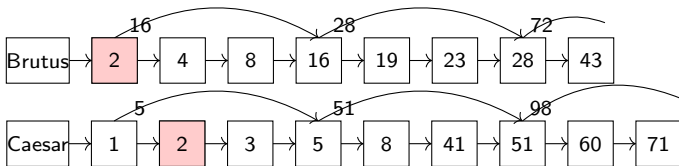


Figure: Postings lists with skip pointers

answer = { }

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer

```

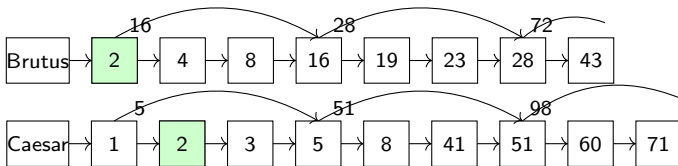


Figure: Postings lists with skip pointers

answer = {2}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

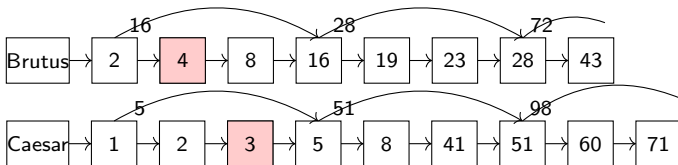


Figure: Postings lists with skip pointers

answer = {2}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

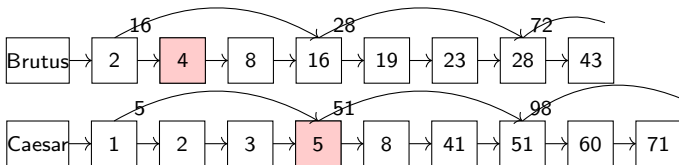


Figure: Postings lists with skip pointers

answer = {2}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

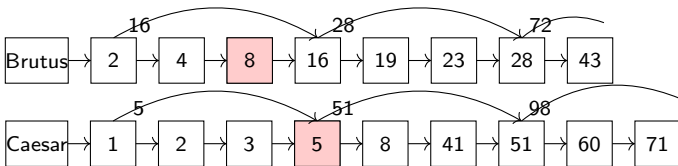


Figure: Postings lists with skip pointers

answer = {2}



```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

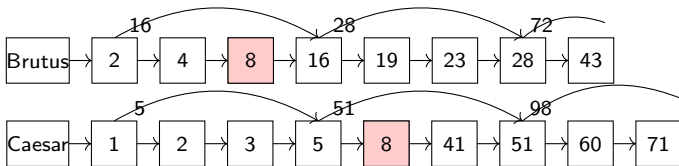


Figure: Postings lists with skip pointers

answer = {2}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer

```

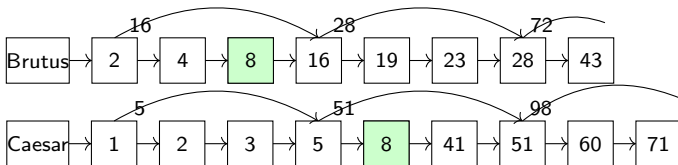


Figure: Postings lists with skip pointers

answer = {2, 8}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

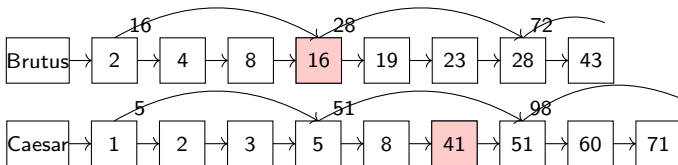


Figure: Postings lists with skip pointers

answer = {2, 8}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

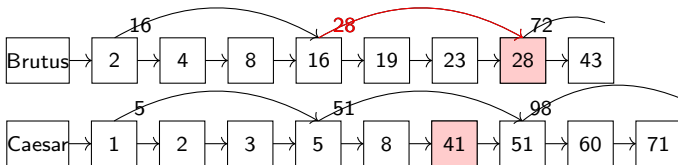


Figure: Postings lists with skip pointers

answer = {2, 8}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

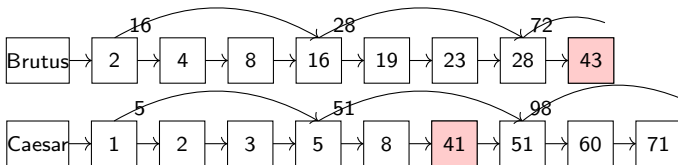


Figure: Postings lists with skip pointers

answer = {2, 8}

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer ← {}
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4  then ADD(answer,  $\text{docID}(p_1)$ )
5      $p_1 \leftarrow \text{next}(p_1)$ 
6      $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8  then if hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9       then while hasSkip( $p_1$ ) and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
10            do  $p_1 \leftarrow \text{skip}(p_1)$ 
11            else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13       then while hasSkip( $p_2$ ) and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
14            do  $p_2 \leftarrow \text{skip}(p_2)$ 
15            else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
    
```

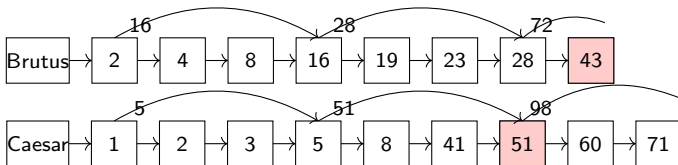


Figure: Postings lists with skip pointers

answer = {2, 8}

## Where to put the skip pointers?

- More skip pointers means shorter skip span, and we are more likely to skip.
- It also means lots of comparisons and lots of spaces to store pointers.
- Fewer pointers means long skip span, and that we are less likely to skip.
- Fewer comparisons and smaller space storing the pointers.
- Use  $\sqrt{P}$  evenly-spaced pointers for a postings list of length  $P$

## Definition

Disjunctive query: a query consisting of terms and logical OR ( $\vee$ ) operator in between

## Problem

*How to process a disjunctive query using inverted indices?*

## Solution

- *locate terms in the dictionary*
- *retrieve their postings*
- *union the postings*



How to union two postings?

How to process a disjunctive query consisting of an arbitrary number of terms?

## Definition

Conjunction normal form (CNF): a Boolean expression is in conjunction normal form if it is a conjunction of some disjunction literals (i.e., it is **AND** of **ORs**)

## Corollary

*Every Boolean expression can be converted to CNF.*

How to convert a Boolean expression consisting of ANDs, ORs, and NOTs into the CNF format?

- 1 Push negations inward by repeatedly applying De Morgan's Law:

- $\neg(P \vee Q) = \neg P \wedge \neg Q$
- $\neg(P \wedge Q) = \neg P \vee \neg Q$
- $\neg\neg P = P$

- 2 Push ORs inward over ANDs by repeatedly apply the distributive law

- $P \vee (Q \wedge R) = (P \vee Q) \wedge (P \vee R)$

If we can process CNF expressions then we can handle arbitrary Boolean queries.


Can we process CNF expressions?

- CNF expressions are conjunctions of disjunctions
- we know how to process a conjunctive query
- we know how to process a disjunctive query

- Boolean information retrieval
- Indexing mechanism
- What is a Boolean query
- How to process Boolean queries
- Term-document Incidence Matrix
- Inverted Index
- Skip-pointers
- Toward processing arbitrary Boolean queries

- Write out a postings merge algorithm for  $x$  OR  $y$ .
- How should the Boolean query  $x$  AND NOT  $y$  be handled?
- *Query optimization* is the process of selecting how to organize the work of answering a query so that the least total amount of work needs to be done by the system. Recommend a query processing order for: (tangerine OR trees) AND (marmalade OR skies) AND (kaleidoscope OR eyes) given the following postings

	<b>term</b>	<b>postings size</b>
	eyes	213312
	kaleidoscope	87009
list sizes:	marmalade	107913
	skies	271658
	tangerine	46653
	trees	316812

 Manning, Christopher D. and Raghavan, Prabhakar and Schütze, Hinrich.

*Introduction to Information Retrieval.*

Cambridge University Press, 2008.

 Wikipedia Article.

*Conjunctive normal form*

[https:](https://en.wikipedia.org/wiki/Conjunctive_normal_form)

[//en.wikipedia.org/wiki/Conjunctive\\_normal\\_form](https://en.wikipedia.org/wiki/Conjunctive_normal_form)