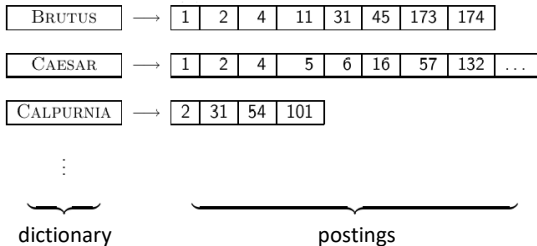


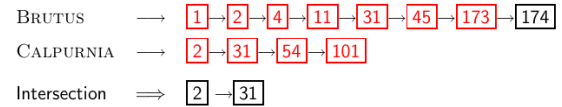
## Inverted Index

For each term  $t$ , we store a list of all documents that contain  $t$ .



1

## Intersecting two posting lists



2

## Westlaw: Example queries

*Information need:* Information on the legal theories involved in preventing the disclosure of trade secrets by employees formerly employed by a competing company  
*Query:* "trade secret" /s disclos! /s prevent /s employe!  
*Information need:* Requirements

for disabled people to be able to access a workplace  
*Query:* discap! /p access! /s work-site work-place (employment /3 place)

*Information need:* Cases about a host's responsibility for drunk guests  
*Query:* host! /p (responsib! liab!) /p (intoxicat! drunk!) /p guest

3

## Does Google use the Boolean model?

- On Google, the default interpretation of a query  $[w_1 w_2 \dots w_n]$  is  $w_1 \text{ AND } w_2 \text{ AND } \dots \text{ AND } w_n$
- Cases where you get hits that do not contain one of the  $w_i$ :
  - anchor text
  - page contains variant of  $w_i$  (morphology, spelling correction, synonym)
  - long queries ( $n$  large)
  - boolean expression generates very few hits
- Simple Boolean vs. Ranking of result set
  - Simple Boolean retrieval returns matching documents in no particular order.
  - Google (and most well designed Boolean engines) rank the result set – they rank good hits (according to some estimator of relevance) higher than bad hits.

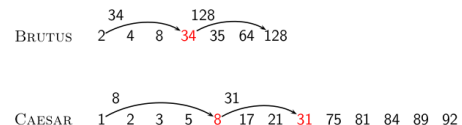
4

## Skip pointers

- Skip pointers allow us to **skip** postings that will not figure in the search results.
- This makes intersecting postings lists more efficient.
- Some postings lists contain several million entries – so efficiency can be an issue even if basic intersection is linear.
- Where do we put skip pointers?
- How do we make sure intersection results are correct?

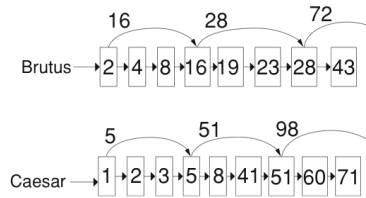
5

## Basic idea



6

## Skip lists: Larger example



7

## 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  $\text{hasSkip}(p_1)$  and ( $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$ )
9  then while  $\text{hasSkip}(p_1)$  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  $\text{hasSkip}(p_2)$  and ( $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$ )
13  then while  $\text{hasSkip}(p_2)$  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

```

8

## Where do we place skips?

- Tradeoff: number of items skipped vs. frequency skip can be taken
- More skips: Each skip pointer skips only a few items, but we can frequently use it.
- Fewer skips: Each skip pointer skips many items, but we can not use it very often.

9

## Where do we place skips? (cont)

- Simple heuristic: for postings list of length  $P$ , use  $\sqrt{P}$  evenly-spaced skip pointers.
- This ignores the distribution of query terms.
- Easy if the index is static; harder in a dynamic environment because of updates.
- How much do skip pointers help?
- They used to help a lot.
- With today's fast CPUs, they don't help that much anymore.

10

## Phrase queries

- We want to answer a query such as [stanford university] – as a phrase.
- Thus *The inventor Stanford Ovshinsky never went to university* should **not** be a match.
- The concept of phrase query has proven easily understood by users.
- About 10% of web queries are phrase queries.
- Consequence for inverted index: it no longer suffices to store docIDs in postings lists.
- Two ways of extending the inverted index:
  - biword index
  - positional index

11

## Biword indexes

- Index every consecutive pair of terms in the text as a phrase.
- For example, *Friends, Romans, Countrymen* would generate two biwords: "friends romans" and "romans countrymen"
- Each of these biwords is now a vocabulary term.
- Two-word phrases can now easily be answered.

12

## Longer phrase queries

- A long phrase like “*stanford university palo alto*” can be represented as the Boolean query “STANFORD UNIVERSITY” AND “UNIVERSITY PALO” AND “PALO ALTO”
- We need to do post-filtering of hits to identify subset that actually contains the 4-word phrase.

13

## Extended biwords

- Parse each document and perform part-of-speech tagging
- Bucket the terms into (say) nouns (N) and articles/prepositions (X)
- Now deem any string of terms of the form  $NX^*N$  to be an *extended biword*
- Examples: catcher in the rye  

$$\begin{array}{ccccc} & N & & X & X & N \\ & \text{king} & \text{of} & \text{Denmark} & & \end{array}$$

$$N \quad X \quad N$$
- Include extended biwords in the term vocabulary
- Queries are processed accordingly

14

## Issues with biword indexes

- Why are biword indexes rarely used?
- False positives, as noted above
- Index blowup due to very large term vocabulary

15

## Positional indexes

- Positional indexes are a more efficient alternative to biword indexes.
- Postings lists in a **nonpositional** index: each posting is just a docID
- Postings lists in a **positional** index: each posting is a docID and a **list of positions**

16

## Positional indexes: Example

Query: “ $to_1 be_2 or_3 not_4 to_5 be_6$ ”  $to$ , 993427:

1: <7, 18, 33, 72, 86, 231>;  
 2: <1, 17, 74, 222, 255>;  
 4: <8, 16, 190, 429, 433>;  
 5: <363, 367>;  
 7: <13, 23, 191>;...

BE, 178239:

1: <17, 25>;  
 4: <17, 191, 291, 430, 434>;  
 5: <14, 19, 101>;... Document 4 is a match!

17

## Proximity search

- We just saw how to use a positional index for phrase searches.
- We can also use it for proximity search.
- For example: employment /4 place
- Find all documents that contain EMPLOYMENT and PLACE within 4 words of each other.
- *Employment agencies that place healthcare workers are seeing growth* is a hit.
- *Employment agencies that have learned to adapt now place healthcare workers* is not a hit.

18

## Proximity search

- Use the positional index
- Simplest algorithm: look at cross-product of positions of (i) EMPLOYMENT in document and (ii) PLACE in document
- Very inefficient for frequent words, especially stop words
- Note that we want to return the actual matching positions, not just a list of documents.
- This is important for dynamic summaries etc.

19

## “Proximity” intersection

```

POSITIONALINTERSECT( $p_1, p_2, k$ )
1  answer ← ( )
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if docID( $p_1$ ) = docID( $p_2$ )
4  then  $I \leftarrow I \cup \{ \}$ 
5       $pp_1 \leftarrow \text{positions}(p_1)$ 
6       $pp_2 \leftarrow \text{positions}(p_2)$ 
7      while  $pp_1 \neq \text{NIL}$ 
8      do while  $pp_2 \neq \text{NIL}$ 
9          do if  $|\text{pos}(pp_1) - \text{pos}(pp_2)| \leq k$ 
10             then ADD( $I, \text{pos}(pp_2)$ )
11             else if  $\text{pos}(pp_2) > \text{pos}(pp_1)$ 
12                 then break
13              $pp_2 \leftarrow \text{next}(pp_2)$ 
14             while  $I \neq \{ \}$  and  $|I[0] - \text{pos}(pp_1)| > k$ 
15                 do DELETE( $I[0]$ )
16             for each  $ps \in I$ 
17                 do ADD(answer, (docID( $p_1$ ), pos( $pp_1$ ), ps))
18              $pp_1 \leftarrow \text{next}(pp_1)$ 
19              $p_1 \leftarrow \text{next}(p_1)$ 
20              $p_2 \leftarrow \text{next}(p_2)$ 
21             if docID( $p_2$ ) < docID( $p_1$ )
22                 then  $p_1 \leftarrow \text{next}(p_1)$ 
23             else  $p_2 \leftarrow \text{next}(p_2)$ 
24  return answer

```

20

## Combination scheme

- Biword indexes and positional indexes can be profitably combined.
- Many biwords are extremely frequent: Michael Jackson, Britney Spears etc
- For these biwords, increased speed compared to positional postings intersection is substantial.
- Combination scheme: Include frequent biwords as vocabulary terms in the index. Do all other phrases by positional intersection.
- Williams et al. (2004) evaluate a more sophisticated mixed indexing scheme. Faster than a positional index, at a cost of 26% more space for index.

21

## “Positional” queries on Google

- For web search engines, positional queries are much more expensive than regular Boolean queries.
- Let's look at the example of phrase queries.
- Why are they more expensive than regular Boolean queries?
- Can you demonstrate on Google that phrase queries are more expensive than Boolean queries?

22