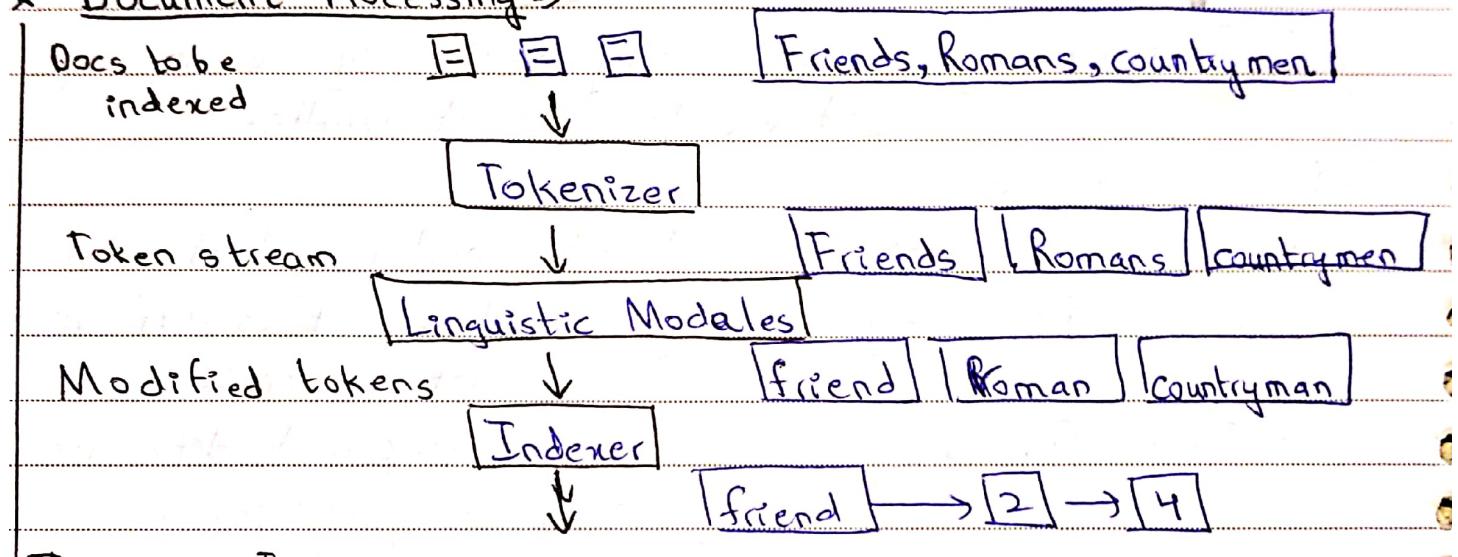


CHAPTER #2

- * Extended Boolean Model => It is an extension to the Boolean Model with additional functionality such as proximity search, phrase search and wildcard search.
- * Word => A delimited string of characters as it appears in text.
- * Term => A "normalized" word (case, morphology, spelling) equivalent of word.
- * Token => An instance of word or a term occurring in a document.
- * Type => same as term usually, equivalence of token.
- * Document Processing =>

Challenges =>

- ↳ format => pdf / word / excel / html

- ↳ language

- ↳ encoding

- * Tokenization => A process through which documents are parsed and a sequence of characters separated, as a feature for document processing.

Issues =>

- ↳ Multi-words => Hewlett-Packard, co-education

- ↳ Numbers => 3/20/91, B-52

- ↳ Languages => dialects → qc, us, >, >, pain

- * Stop words => With a stop words list, we can exclude the most common words entirely. [the, a, is]
 - ↳ Trend: is away from ignoring them
 - ↳ Good compression techniques => space for including stopwords in a system is very well.
 - ↳ Good query optimization => pay little query time for including stop words.

↳ Why they are important?

↳ Phrase queries => "King of Denmark"

↳ Song titles => "Let it be", "To be or not to be"

↳ Relational queries => "flights to London"

- * Normalization => Normalize words in indexed text as well as query words into the same basic form.

↳ Some strategies

↳ Deleting periods => U.S.A → USA

↳ Deleting hyphens => co-education → co.education

↳ Accents => résumé → resume

- * Case folding => Reduce all letters to lowercase.

↳ General Motors → general motors

- * Asymmetric Expansion => Expand the search query to include related terms. It increases recall but reduces precision.

↳ window → window, windows

↳ windows → Windows, windows, window

- * Equivalence Classing => Divide the search space into a smaller number of equivalent classes, each of which represent a group of related terms.

↳ {run, running, ran, runner} → single class

* Thesauri => A lexical resource that provides a collection of related words, including synonyms, antonyms and related terms.

↳ How to handle synonyms?

↳ hand constructed equivalence classes → car = automobile

↳ rewrite to form equivalence classes terms

↳ if a doc has automobile, index it under car too.

↳ expand a query

↳ if a query has automobile, search under car too.

* Soundex => A phonetic algorithm that is used to index words by their sound. The goal is to group similar sound words [homonyms].

↳ Why soundex?

↳ It can handle spelling mistakes since it just checks how words sound.

* Lemmatization => Implies doing "proper" reduction of inflectional / variant forms to dictionary head word base forms.

↳ Examples

↳ am, are, is → be

↳ car, cars, car's, car's' → car

↳ Definition

↳ algorithmic process of determining lemma from a given word.

↳ It may involve

↳ understanding context and determining the parts of speech.

↳ 'the walk' → also written as {walk, walked, walks}

↳ the base form walk is called lemma.

↳ Hard to implement

↳ in case of new languages

* Stemming \Rightarrow Stemming operates on a single word without knowledge of context.

\hookrightarrow Easier to implement

\hookrightarrow run faster but with reduced accuracy.

\hookrightarrow Porter Stemmer

\hookrightarrow word is cleaned up in the initialization phase, one prefix trimming phase, then five suffix trims.

\hookrightarrow Initialization

\hookrightarrow convert to lowercase and keep only digits & letters.

\hookrightarrow F-16 \rightarrow f16

\hookrightarrow Remove prefixes

\hookrightarrow {kilobyte, megabyte, nanobyte} \rightarrow byte

\hookrightarrow Suffix trim 1

\hookrightarrow remove "es" from words ending in "sses" or "ies"

\hookrightarrow passes \rightarrow pass, cries \rightarrow cri

\hookrightarrow remove "s" from words where second last letter is not "s"

\hookrightarrow runs \rightarrow run, fuss \rightarrow fuss

\hookrightarrow remove "ed" if word has a vowel and ends with "eed".

\hookrightarrow agreed \rightarrow agre, freed \rightarrow freed.

\hookrightarrow replace trailing "y" with "i" if word has a vowel

\hookrightarrow satisfy \rightarrow satisfi, fly \rightarrow fli

\hookrightarrow Suffix trim 2

\hookrightarrow replace any suffix from left list to right list.

\hookrightarrow tional \rightarrow tion, ization \rightarrow ize, fulness \rightarrow ful

\hookrightarrow Suffix trim 3

\hookrightarrow replace any suffix from left to right list.

\hookrightarrow icate \rightarrow ic, ative \rightarrow -

\hookrightarrow Suffix trim 4

\hookrightarrow remove al, ance, ence, etc..

\hookrightarrow Suffix trim 5 [hinge \rightarrow hing]

\hookrightarrow remove trailing 'e' if word does not end in a vowel

* Stemming	V/S	Lemmatization
1) Heuristic, rule based approach, generally fast, single term.		1) Rigor process that uses a dictionary and context to determine lemma, a slow approach.
2) generates unreadable tokens.		2) generates readable lexeme.
3) too aggressive, sacrificing precision for recall.		3) sacrifices recall for precision.
* Morphology \Rightarrow Field of linguistics focused on the study of forms and formations of words in a language.		
↳ Inflectional \Rightarrow adding a suffix to a word that does not change its grammatical category.		
↳ work \rightarrow worked, working		
↳ Derivational \Rightarrow adding a suffix to a word that changes its grammatical category.		
↳ nation (noun) \Rightarrow national (adj) \Rightarrow nationalize (verb)		
↳ Implementation issues in inverted Inverted Index		
↳ Lists \Rightarrow Lists can suffer from slow performance when dealing with complex morphological forms due to linear search. The memory reqs for storing frequency information for each term can become high when dealing with morphologically rich terms.		
↳ Hashmaps \Rightarrow The hash function must be able to deal with morphisms without causing collisions. The scalability of the hash function is also a concern due to increasing data.		
↳ Trees \Rightarrow Maintaining a balanced tree can be challenging and requires a deep understanding of algorithms and data structures.		

- ↳ Implementation issues with skip lists
- ↳ Handling complex morphs => Morphologically rich data is tough to sort - Inflections and derivations can make the order tough to determine, thus lexicographical sort is harder.
- ↳ Scalability => They are designed to be scalable but the size of the data can increase overtime making maintenance tough.
- ↳ Memory management => They require a large amount of memory to store each frequency for morphologically rich terms since the memory requirements for such terms becomes prohibitively high.

* Phrase queries => The multiword compounds or phrase queries make sense together e.g.: "Stanford University". Such words cannot be broken down separately and queries due to false positives retrieved. So, the following solutions are implemented:

↳ Bi-word Indexes => The first approach is to index every consecutive pair of terms in the text as a phrase. For example, the text "Friends, Romans, Countrymen" will generate the biwords "friends romans", "romans countrymen". The two word phrase query-processing is now immediate.

↳ Longer Phrase queries => In cases of longer phrases like "stanford university pao alto", the following query will be broken into "stanford university", "university pao", "pao alto". Without the actual documents, we cannot verify if the actual phrase was retrieved, hence there will be a lot of false positives.

phrase index

↳ Extended bi-words ⇒ We firstly parse the indexed text and part-of-speech tagging (POST). [a process that involves assigning a part of speech to each word]. Then, group the words into nouns (N) and function words (articles/prepositions) (X). Now, each string of the form $N X^* N$ is an extended bi-word. For example: Coins are in pocket, the

$N \ X \ X \ N$

lookup and dictionary word for the text will be "coin pocket".

↳ Issues ⇒

↳ False +ves.

↳ Index blowup due to extremely huge dictionary.

↳ Positional Indexes ⇒ In the posting lists, store the position of appearance of each term for a document as well.

$\langle \text{term}, \text{number of docs}; \text{docNumber}; \text{position}_1, p_2; \rangle$

For phrase queries, we would have to use a merge algorithm recursively while dealing with term proximity as well.

The following algorithm uses positional indexes to find phrases within proximity.

POSITIONALINTERSECT (p_1, p_2, k):

- 1 answer ← ()
- 2 while p_1 and p_2 : matching IDs
- 3 do if $\text{docID}(p_1) = \text{docID}(p_2)$
- 4 l ← ()
- 5 $pp_1 \leftarrow \text{positions}(p_1)$] get term
- 6 $pp_2 \leftarrow \text{positions}(p_2)$] positions

b be 2 ← 5, 6 ①^{to}
1 2 3 4 5 6 7 8 4

Dated: _____

7 while pp1:
8 do while pp2:
9 do if $|pos(pp1) - pos(pp2)| \leq k$ proximity
10 add (2, pos(pp2)) distance
11 else if $pos(pp2) > pos(pp1)$
12 break since sorted, no need to
13 go ahead
14 some sort [pp2 ← next(pp2) ← get next position]
15 of reckoning [while $|l| \neq 1$ and $|l[0] - pos(pp1)| > k$
16 for each pse l
17 add (answer, < docID(p1), pos(pp1), ps>)
18 pp1 ← next(pp1)
19 p1 ← next(p1)] get next
20 p2 ← next(p2)] posting list
21 else if docID(p1) < docID(p2)
22 p1 ← next(p1)
23 else
24 p2 ← next(p2)
25 return answer.

* Proximity Queries => They also use the above intersection algorithm to search words within proximity of $\backslash k$ words. e.g.: STATUE /3 FEDERAL /2 TORT. Clearly positional indexes are required for such queries.

* Positional Index Size => Although the positional index expands the posting storage substantially, it is still standardized due to its power of proximity query searching.
↳ Rules of thumb => { holds for all English-like languages }

↳ A positional index is 2-4 as large as non-positional index.
↳ It is 35-50% of the volume of original text.

* Combination Schemes \Rightarrow If users query commonly on particular phrases, such as "Michael Jackson", it is inefficient to keep merging positional posting lists.

A combination strategy uses binword index for certain queries and positional index for the rest. Good queries are to be involved as phrase indexes such that the bi-words together are known to be common based on percent querying behavior.

\hookrightarrow Williams' Strategy \Rightarrow He combined the two schemes such that when a query is made, the phrase indexes are used to identify potential matches then the positional index is used to rank those matches based on closeness. With this strategy, he achieved the precision of positional indexing while being able to handle phrases. It required 2.61 times more space but executed a simple web query in $1/4$ time of positional index.

* Query Processing Skip Lists \Rightarrow We can use skip lists for faster query processing.

\hookrightarrow Tradeoffs

\hookrightarrow More skips \rightarrow shorter skip spans \Rightarrow more likely to skip but lots of comparison to skip pointers.

\hookrightarrow Fewer skips \rightarrow long skip spans \Rightarrow few successful skips

\hookrightarrow Placing skips

\hookrightarrow for postings of length L, use \sqrt{L} evenly-spaced skip pointers.

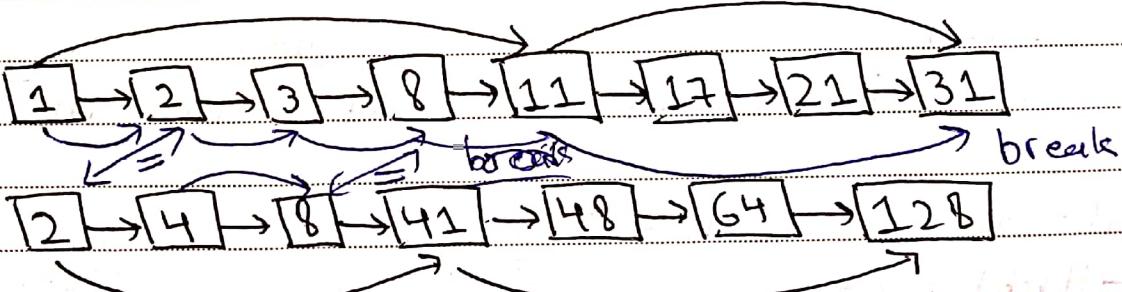
\hookrightarrow Hard to maintain if L is dynamic.

Dated: _____



INTERSECT WITH SKIPS (p_1, p_2):

```
1 answer ← ()  
2 while  $p_1$  and  $p_2$ :  
3     if docID( $p_1$ ) = docID( $p_2$ ):  
4         add(answer, docID( $p_1$ )) ] if equal,  
5          $p_1 \leftarrow \text{next}(p_1)$  ] match found  
6          $p_2 \leftarrow \text{next}(p_2)$   
7     else if docID( $p_1$ ) < docID( $p_2$ ):  
8         while skip } if hasSkip( $p_1$ ) and docID(skip( $p_1$ )) ≤ docID( $p_2$ )  
9             is possible } while hasSkip( $p_1$ ) and docID(skip( $p_1$ )) ≤ docID( $p_2$ )  
10            keep skipping }  $p_1 \leftarrow \text{skip}(p_1)$   
11     else  
12    else go one by one  $p_1 \leftarrow \text{next}(p_1)$   
13    { else if hasSkip( $p_2$ ) and docID(skip( $p_2$ )) ≤ docID( $p_1$ )  
14        same } while hasSkip( $p_2$ ) and docID(skip( $p_2$ )) ≤ docID( $p_1$ )  
15        as }  $p_2 \leftarrow \text{skip}(p_2)$   
16        above } else  
17         $p_2 \leftarrow \text{next}(p_2)$   
18    return answer
```



→ What is the document processing pipeline for IR system?

It refers to a series of computational steps that are performed on a set of textual documents to transform them into a format that is suitable for indexing and searching.

The stages involve:

1) Document Acquisition ⇒ System collects raw text docs,

2) Text preprocessing ⇒ Cleaning and normalizing docs by removing irrelevant documents such as HTML tags, punctuation, stop words. It also involves stemming and lemmatization.

3) Indexing ⇒ The preprocessed text is converted into a form that can be efficiently stored or searched.

4) Query Processing ⇒ This stage involves processing user queries, typically in the form of a set of keywords and transforming them into a form that can be matched against the index. This may involve similar preprocessing and indexing steps for the query.

5) Ranking & Retrieval ⇒ The search engine retrieves a ranked list of relevant documents based on the query and similarity to the query.