



Inverted Index

Dr. Subrat Kumar Nayak

Associate Professor

Dept. of CSE, ITER, SOADU

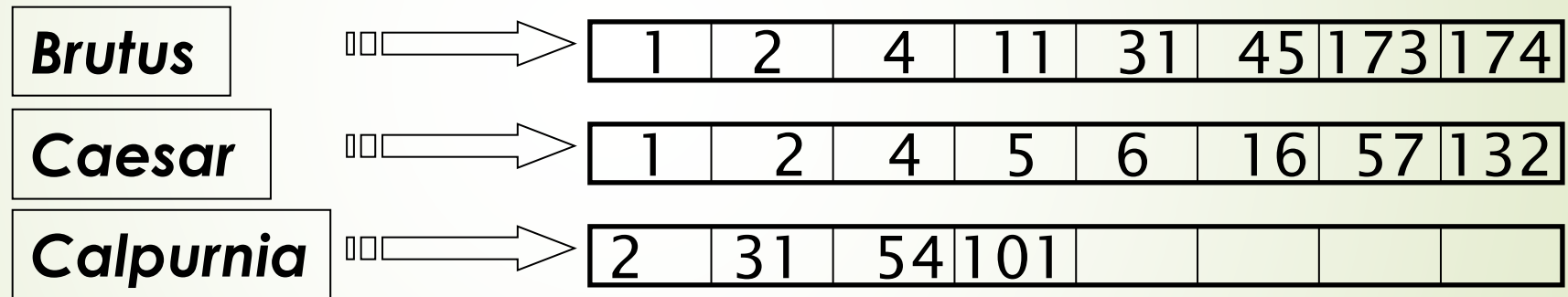
Inverted Index

- This is the first major concepts of information retrieval.
- The name is actually redundant: **an index** always maps back from terms to the parts of a document where they occur.
- Inverted Index, sometimes coined as **Inverted file**.
- **Inverted Index is used** keep a *dictionary* of terms. Then **for each term**, we have a list that records which documents the term occurs in.
- Each **item in the list**, which records that a term appeared in a document is conventionally called a *posting*.
- The list is then called a *postings list* (or **inverted list**), and all the postings lists taken together are referred to as the *postings*.

dictionary term is used for the data structure and *vocabulary* for the set of terms

Inverted Index

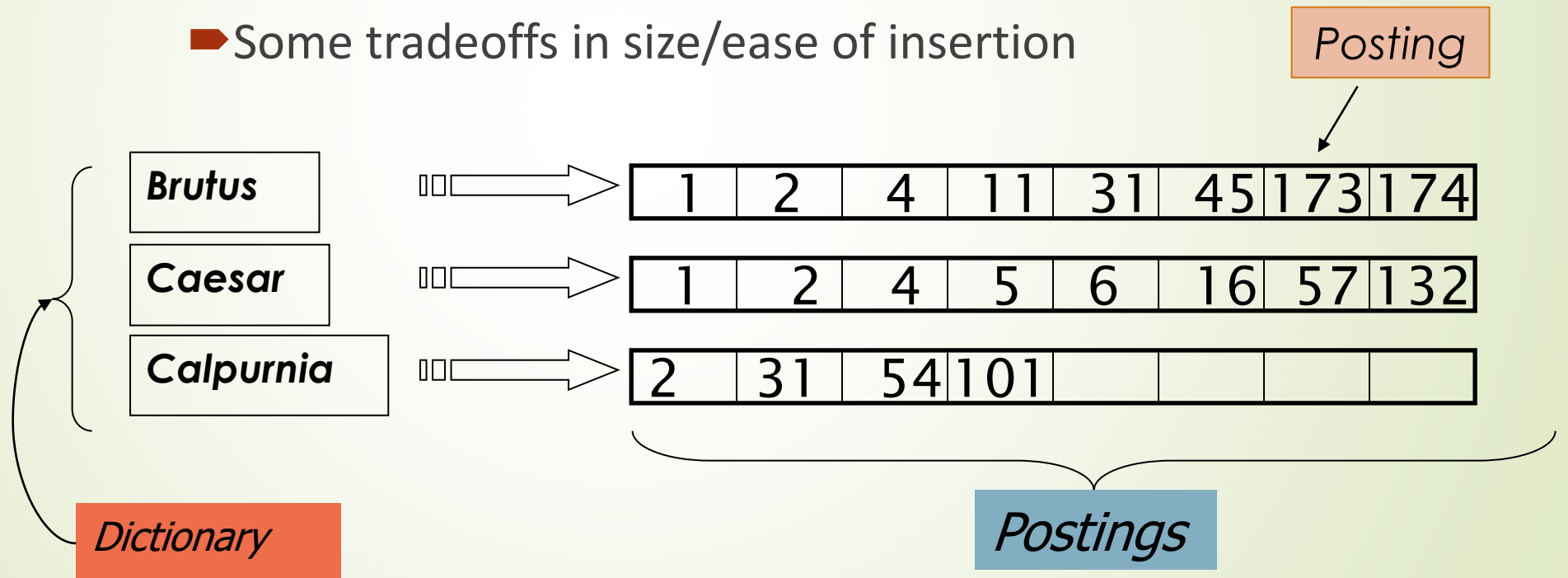
- For each term t , we must store a list of all documents that contain t .
 - Identify each doc by a **docID**, a document serial number
- Can we use fixed-size arrays for this?



What happens if the word **Caesar** is added to document 14?

Inverted Index

- We need variable-size **postings lists**
 - On disk, a continuous run of postings is normal and best
 - In memory, can use linked lists or variable length arrays
 - Some tradeoffs in size/ease of insertion



Sorted by docID (more later on why).

Inverted Index (Construction)

Documents to
be indexed



Friends, Romans, countrymen.

⋮

Tokenizer

Token stream

Friends

Romans

Countrymen

Linguistic modules

Modified tokens

friend

roman

countryman

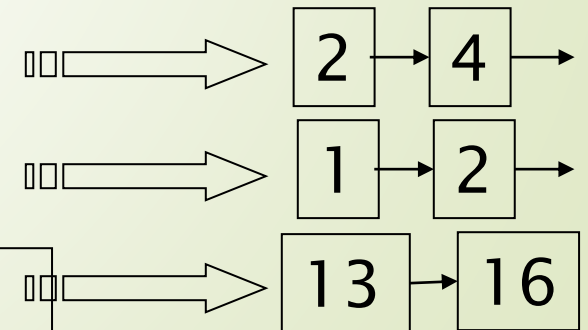
Indexer

Inverted index

friend

roman

countryman



Initial stages of text processing

- Tokenization
 - Cut character sequence into word tokens
 - Deal with “*John’s*”, *a state-of-the-art solution*
- Normalization
 - Map text and query term to same form
 - You want *U.S.A.* and *USA* to match
- Stemming
 - We may wish different forms of a root to match
 - *authorize, authorization*
- Stop words
 - We may omit very common words (or not)
 - *the, a, to, of*

Indexer steps: Token sequence

- Sequence of (Modified token, Document ID) pairs.

Doc 1

I did enact Julius
Caesar I was killed
i' the Capitol;
Brutus killed me.

Doc 2

So let it be with
Caesar. The noble
Brutus hath told you
Caesar was ambitious

Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2

Indexer steps: Sort

- Sort by terms
 - At least conceptually
 - And then docID

Core indexing step

Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2



Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2

Indexer steps: Dictionary & Postings

- Multiple term entries in a **single document** are merged.
- Split into Dictionary and Postings
- Doc. frequency information is added.

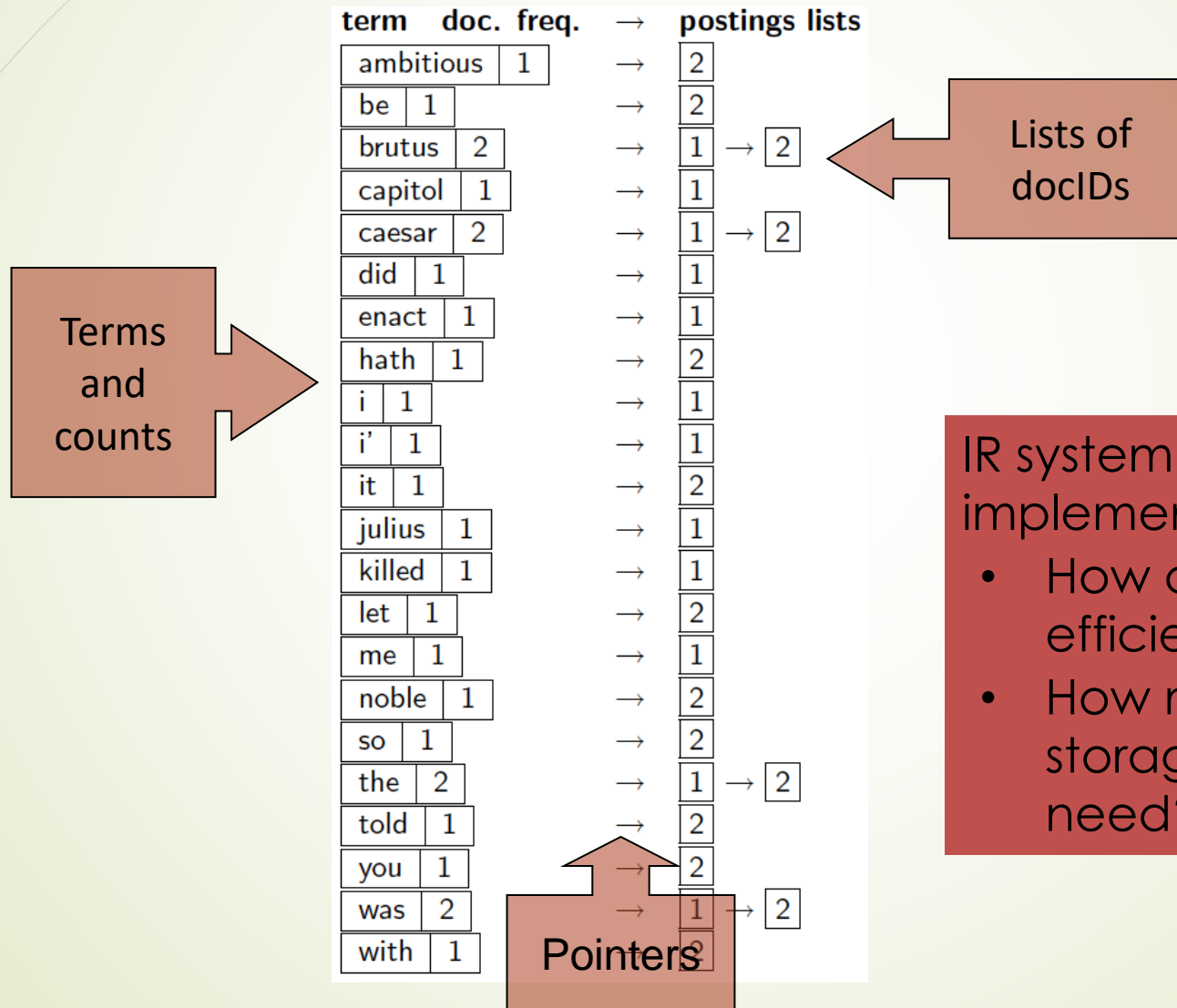
Why frequency?
Will discuss later.

Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2



term	doc. freq.	→	postings lists
ambitious	1	→	2
be	1	→	2
brutus	2	→	1 → 2
capitol	1	→	1
caesar	2	→	1 → 2
did	1	→	1
enact	1	→	1
hath	1	→	2
i	1	→	1
i'	1	→	1
it	1	→	2
julius	1	→	1
killed	1	→	1
let	1	→	2
me	1	→	1
noble	1	→	2
so	1	→	2
the	2	→	1 → 2
told	1	→	2
you	1	→	2
was	2	→	1 → 2
with	1	→	2

Where do we pay in storage?

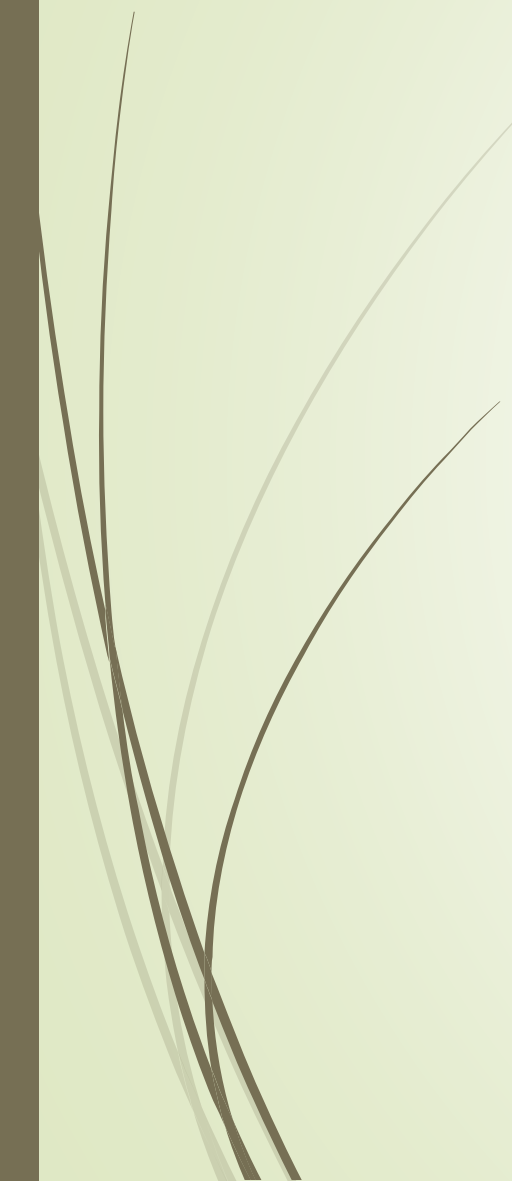


IR system implementation

- How do we index efficiently?
- How much storage do we need?

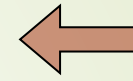


Inverted Index

- Inverted index works much better than the Boolean retrieval method.
 - Sorting based inverted indexing is more efficient than the inverted indexing method since least work needs to be done.
- 

Query processing with an inverted index

- How do we process a query?
 - Later – what kinds of queries can we process?



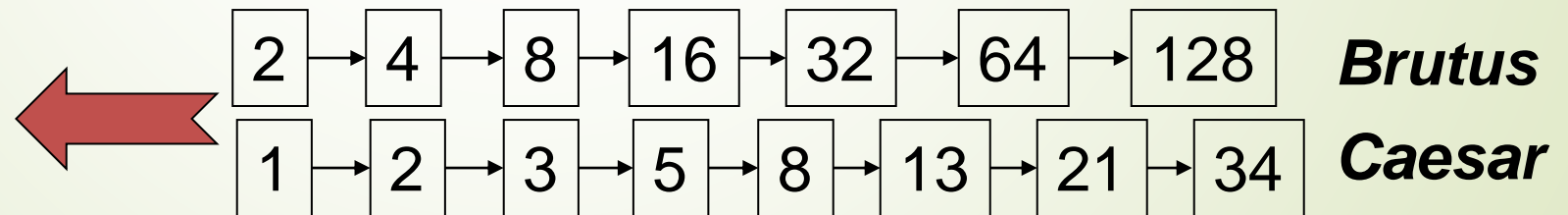
Our focus

Query processing: AND

- Consider processing the query:

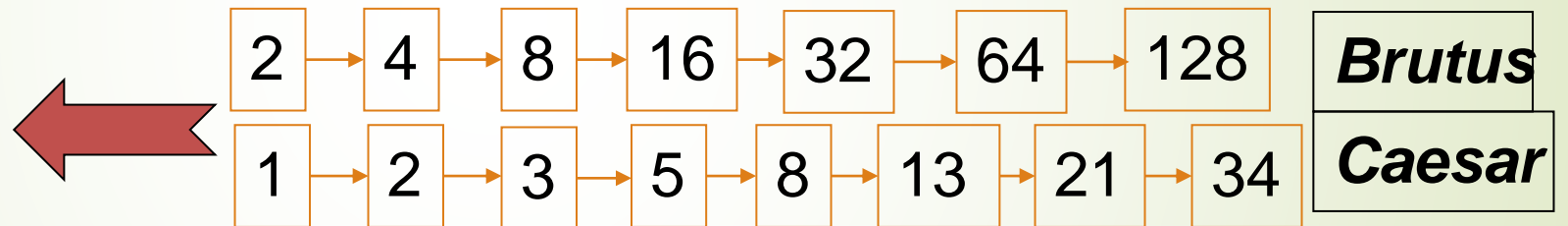
Brutus AND Caesar

- Locate **Brutus** in the Dictionary;
 - Retrieve its postings.
- Locate **Caesar** in the Dictionary;
 - Retrieve its postings.
- “Merge” the two postings (intersect the document sets):



The merge

- Walk through the two postings simultaneously, in time linear in the total number of postings entries



If the list lengths are x and y , the merge takes $O(x+y)$ operations.

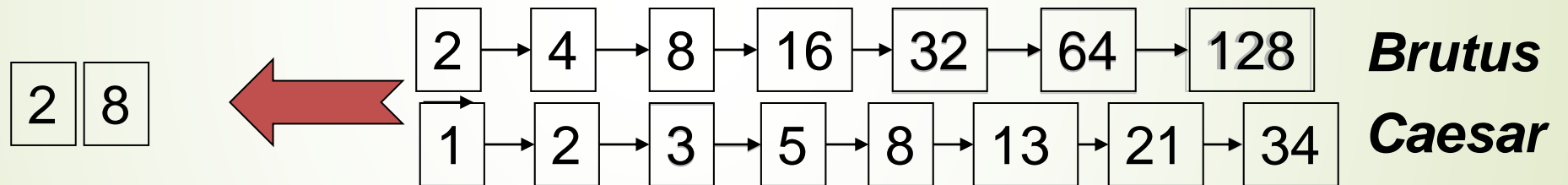
Crucial: postings sorted by docID.

Intersecting two postings lists (a “merge” algorithm)

```
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  $answer$ 
```


The merge

- Walk through the two postings simultaneously, in time linear in the total number of postings entries



If the list lengths are x and y , the merge takes $O(x+y)$ operations.

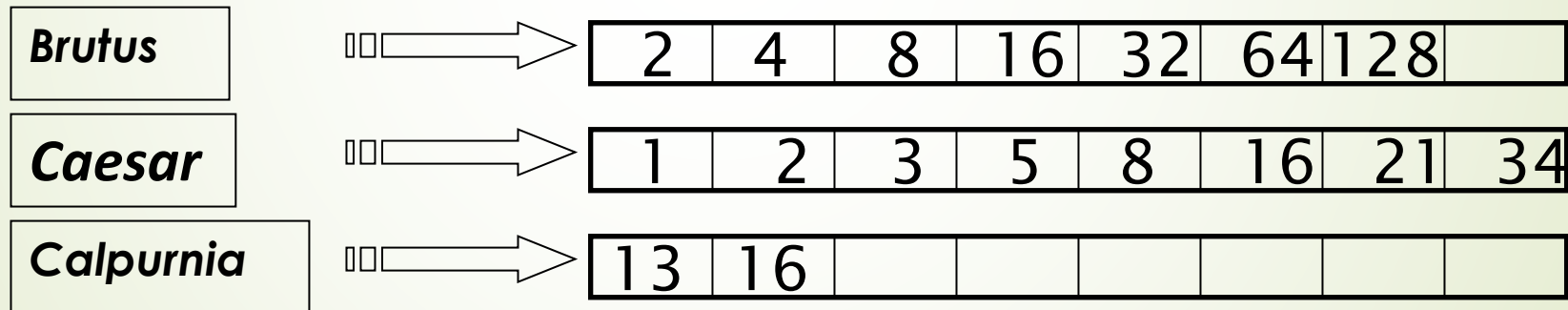
Crucial: postings sorted by docID.

Boolean queries: Exact match

- The **Boolean retrieval model** is being able to ask a query that is a Boolean expression:
 - Boolean Queries are queries using *AND*, *OR* and *NOT* to join query terms
 - Views each document as a set of words
 - Is precise: document matches condition or not.
 - Perhaps the simplest model to build an IR system on
- Primary commercial retrieval tool for 3 decades.
- Many search systems you still use are Boolean:
 - Email, library catalog, macOS Spotlight

Query optimization

- ▶ *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.
- ▶ What is the best order for query processing?
- ▶ Consider a query that is an *AND* of n terms.
- ▶ For each of the n terms, get its postings, then *AND* them together.



Query: **Brutus AND Calpurnia AND Caesar**

Query optimization example

- ▶ Process in order of increasing freq:
 - ▶ *start with smallest set, then keep cutting further.*

This is why we kept
document freq. in dictionary

Brutus	→	2	4	8	16	32	64	128	
Caesar	→	1	2	3	5	8	16	21	34
Calpurnia	→	13	16						

Execute the query as (***Calpurnia AND Brutus***) ***AND Caesar***.

More general optimization

e.g., (*madding OR crowd*) AND (*ignoble OR strife*)

- Get doc. freq.'s for all terms.
- Estimate the size of each *OR* by the sum of its doc. freq.'s (conservative).
- Process in increasing order of *OR* sizes.

Algorithm to Intersect n 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** $result$

