

INFORMATION RETRIEVAL

Dr. Reda M. Hussien

Term-document incidence matrices

Unstructured data in 1620

- Which plays of Shakespeare contain the words Brutus AND Caesar but NOT Calpurnia?
- One could grep all of Shakespeare's plays for Brutus and Caesar, then strip out lines containing Calpurnia?
- Why is that not the answer?
 - Slow (for large corpora)
 - NOT Calpurnia is non-trivial
 - Other operations (e.g., find the word Romans near countrymen) not feasible
 - Ranked retrieval (best documents to return)
 - Later lectures

Term-document incidence matrices

| | Antony 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 |

Brutus AND Caesar BUT NOT Calpurnia

1 if play contains word, 0 otherwise



Incidence vectors

- So we have a 0/1 vector for each term.
- To answer query: take the vectors for Brutus, Caesar and Calpurnia (complemented)
→ bitwise AND.

- 110100 AND
- 110111 AND
- 101111 =
- 100100

| | Antony 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 |

Answers to query

- Antony and Cleopatra, Act III, Scene ii

Agrippa [Aside to DOMITIUS ENOBARBUS]: Why, Enobarbus,

When Antony found Julius **Caesar** dead,

He cried almost to roaring; and he wept

When at Philippi he found **Brutus** slain.

- Hamlet, Act III, Scene ii

Lord Polonius: I did enact Julius **Caesar** I was killed i' the

Capitol; **Brutus** killed me.

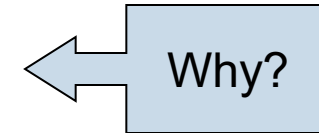


Bigger collections

- Consider $N = 1$ million documents, each with about 1000 words.
- Avg 6 bytes/word including spaces/punctuation
 - 6GB of data in the documents.
- Say there are $M = 500K$ distinct terms among these.

Can't build the matrix

- 500K x 1M matrix has half-a-trillion 0's and 1's.
- But it has no more than one billion 1's.
 - matrix is extremely sparse.
- What's a better representation?
 - We only record the 1 positions.



Inverted index

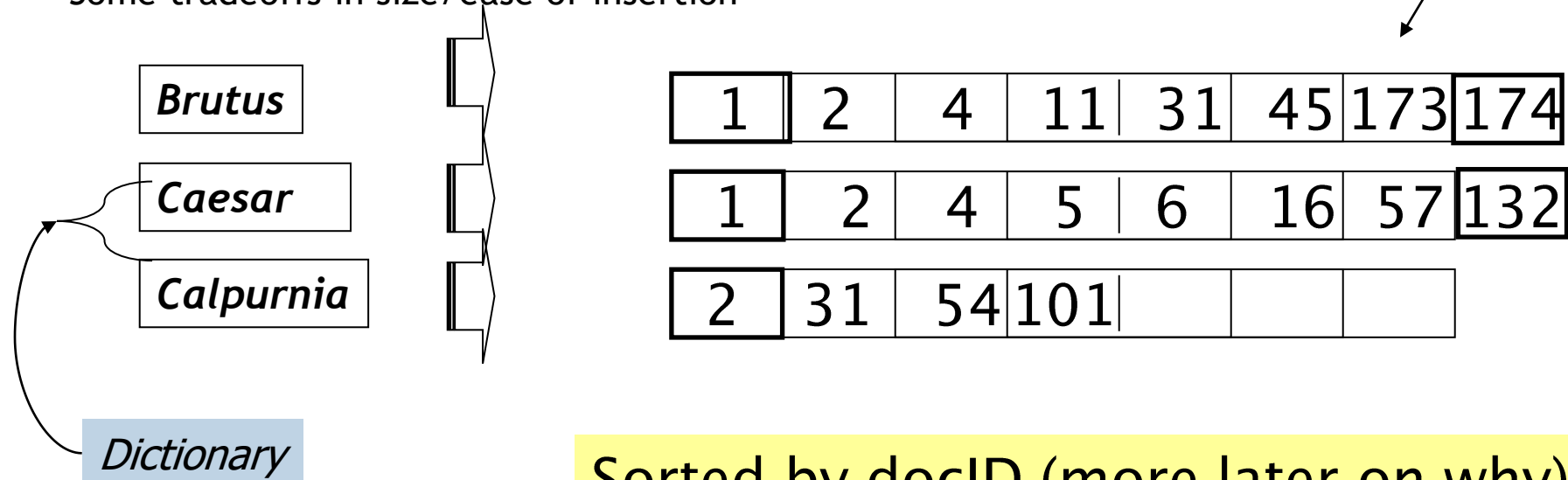
- For each term t , we must store a list of all documents that contain t .
 - Identify each doc by a docID, a document serial number
- Can we use fixed-size arrays for this?

| | | | | | | | | | |
|-------------------------|--|---|----|----|-----|----|----|-----|-----|
| <i>Brutus</i> | | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
| <i>Caesar</i> | | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 |
| <i>Calpurnia</i> | | 2 | 31 | 54 | 101 | | | | |

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

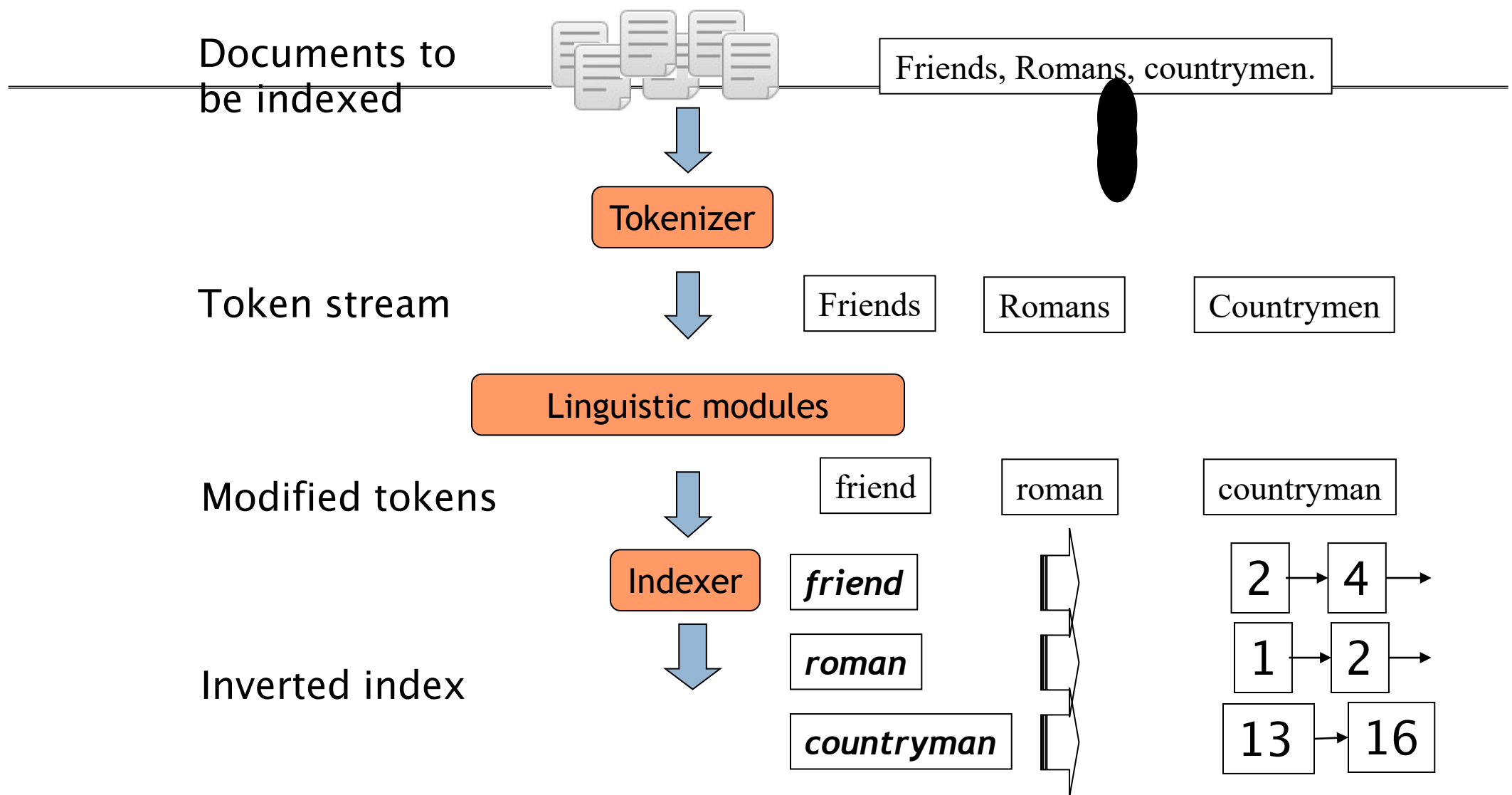
Inverted index

- 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



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

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

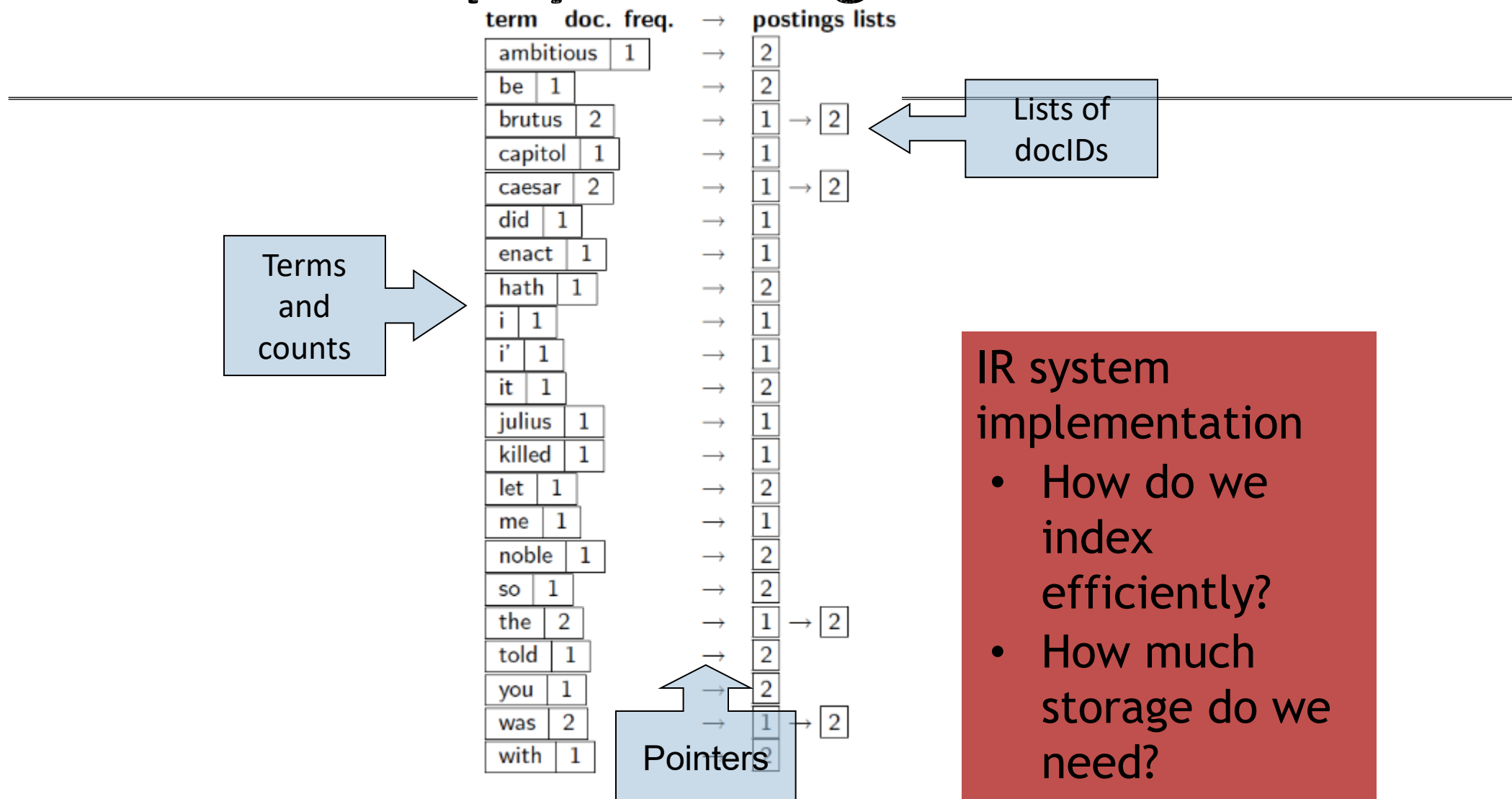
| Term | docID |
|-----------|-------|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 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?