

Boolean Retrieval.

- The query is expressed as a Boolean expression

- terms are the operands.

- standard Boolean operators
AND, OR, NOT.

Corpus: collection of plays written by Shakespeare.

Document: an individual play.

~~Query~~ Query: Boolean expression having terms connected by logical operators.

- Which plays of Shakespeare contain words Brutus AND Caesar but NOT Calpurnia

Design of Information

Retrieval Systems.

Components/terms related to the design of IR Systems:

- Documents (docId)
- Corpus (collection of documents)
- **User** has information need
↳ query string
- "Term" ← basic unit of information
- Relevance

algorithmic challenges.

Known world
documents

unknown world
query.

relevant?

Brute force:

grep Brutus & Caesar

↓
result → all lines that
do not have
California.

Challenges:

— text volume is too large.

Brutus AND Caesar	NOT California.
-------------------	-----------------

↓ one time confirmation

↓ more exhaustive

"near"

words

Roman/Romans

near countrymen

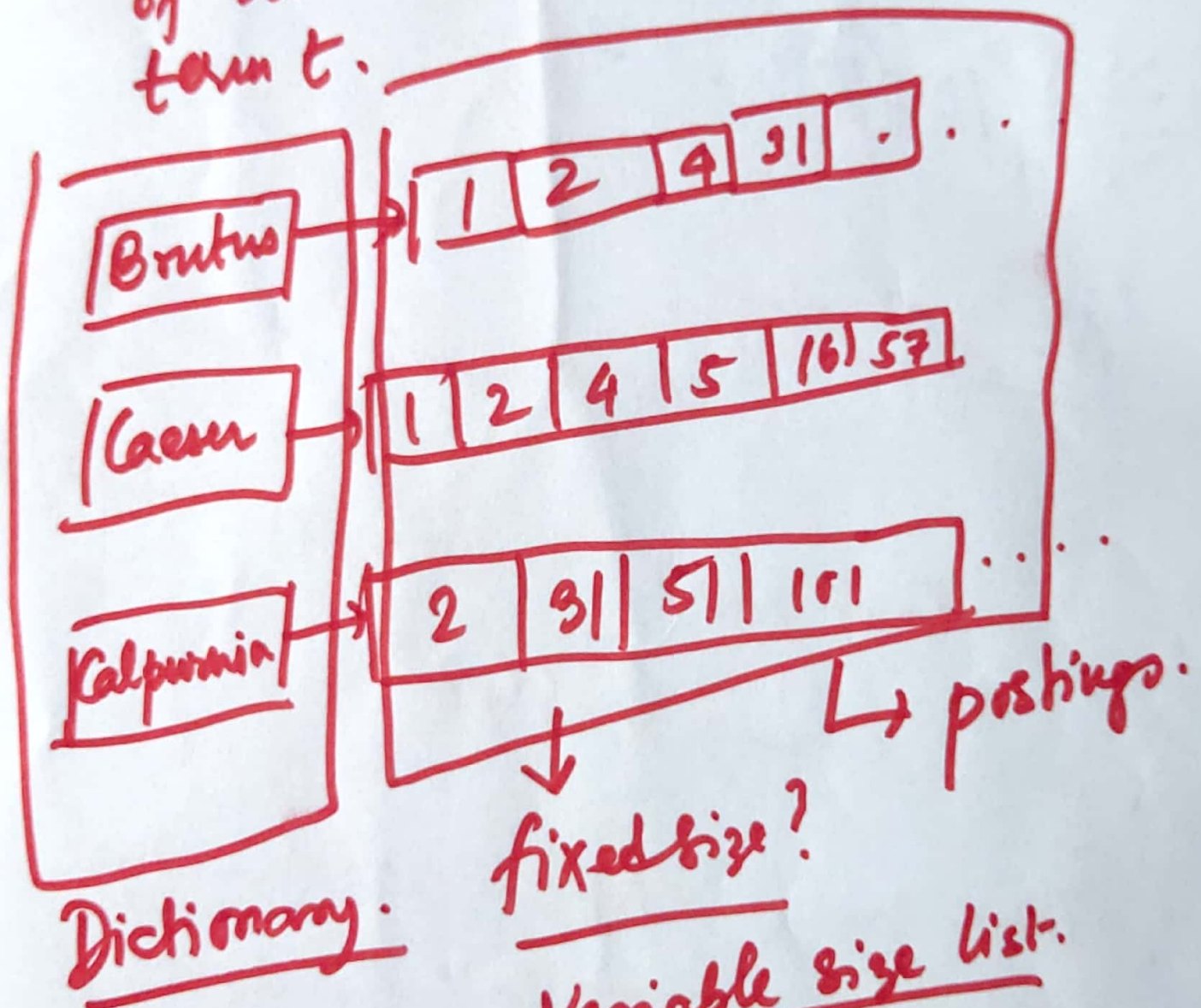
Term-document incidence
matrix

	Ant	Je	T Temp	Hamlet	Phello	Mac
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
Mercury	1	0	1	1	1	1
Horser	1	0	1	1	1	0

1000 x 14 < 18.

Inverted Index.

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



Variable size list.

↓
posting list

Boutus: 110100 AND

Caesar: 110111 AND

Calpurnia: 101111

\Rightarrow 100100. \leftarrow

$N = 1M$ docs.

\hookrightarrow 1000 words.

Each word \rightarrow 6 bytes. (on avg)

$\hookrightarrow \sim 6GB$.

There are $M = 500K$ unique terms across all the docs.

$500K \times 1M \sim \frac{1}{2}$ trillion.

How many 1's \rightarrow 1 billion