

Tokenization: cutting/splitting char into
"word token" ← terms.

Language Model: a) lower casing
b) stemming ←

U.S.A.

USA

Same identity (entity)

Normalization:

Sequence (Modified token, doc id).

- sort by the terms as they occur in the document
- Sort by the doc id.

Index Built!

Process Query.

Boulus AND Caesar.

Locate Brulius in dict.

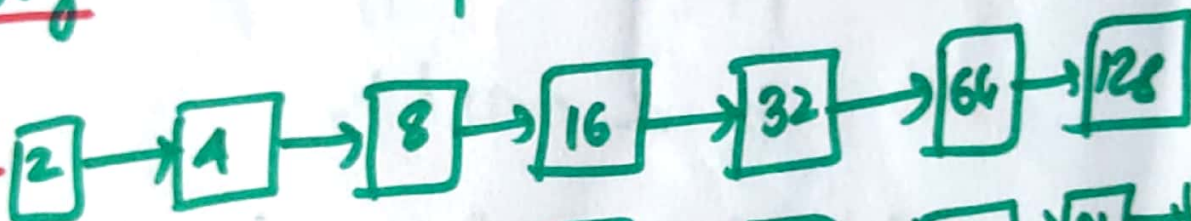
Locate Caesar in dict.

retrieve
postings

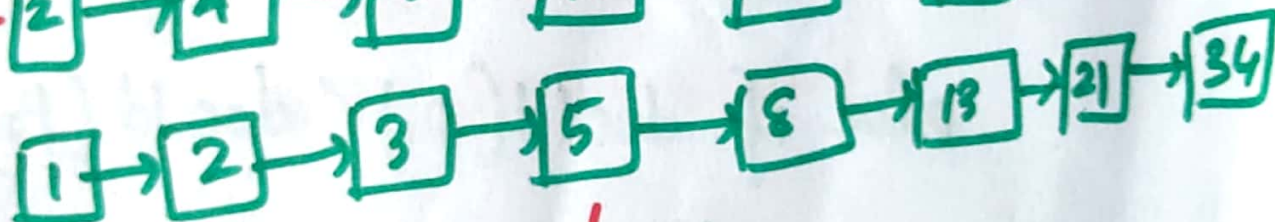
retrieve postings.

"Merge" the two postings

Brulius



Caesar



doc ids are sorted in the
posting list.

$O(x+y)$.

Intersect (P_1, P_2)

answer $\leftarrow \langle \rangle$

while $P_1 \neq \text{NIL}$ & $P_2 \neq \text{NIL}$

do if $\text{docId}(P_1) = \text{docId}(P_2)$

then $\text{ADD}(\text{answer}, \text{docId}(P_1))$

$P_1 \leftarrow \text{next } P_1$

$P_2 \leftarrow \text{next } P_2$

else if $\text{docId}(P_1) < \text{docId}(P_2)$

then $P_1 \leftarrow \text{next } P_1$

else $P_2 \leftarrow \text{next } P_2$

return answer.

Query optimization

What is the best order in which we process our query.

query \leftarrow simple AND of some n query terms.

(document frequency or the term) ~~term~~ frequency \rightarrow indicator of the size of the posting ~~list~~ list. of the term across documents

two terms that have the smallest doc frequency.

OR queries?

(madding OR crowd) AND (ignoble AND strife)

Estimate the size of the OR?

worst case — sum of the doc frequencies.

PROCEEDS IN INCREASING ORDER OF OR SIZES

How popular do you think
is Boolean retrieval?

- Not very popular.
- somewhat popular.
- used to be ^{at the peak of} popularity

Legal IR system:

WestLaw

www.westlaw.com.

700K users.

- What is the statute of limitations
in cases involving the federal
tort claims act?

↓
LIMIT! / 3 STATUTE ACTION / S
FEDERAL / 2 TORT / 3 CLAIM.

/3 → word windows (within 3 words)
/5 → within the sentence

Q → Google's default interpretation
of the query $[w_1 w_2 w_3 w_4] \dots$

Default interpretation.

w_1 AND w_2 AND w_3 AND $w_4 \dots$

When does this not hold?

- (i) page contains a variant of w_i
- (ii) presence of w_i in anchor text.

Order or ambiguity
of the query.

→ SKY AND blue | ← color.

→ blue AND SKY | ← video

→ common AND wealth | — Games

→ wealth AND common | → "common sense is your wealth"