

WILDCARD QUERIES

WILDCARD QUERIES SUCH AS

mon^* OR $*mon$ ARE EASY TO
DO WITH PREFIX SEARCH:

- α^* = find all documents containing words starting with α

SIMPLE PREFIX SEARCH

- $*\beta$ = find all documents containing words ending in β

PREFIX SEARCH, BUT YOU HAVE TO CONSTRUCT
A DATA STRUCTURE CONTAINING THE REVERSE
OF EACH TERM

BUT WHAT IF WE SEARCH FOR

$\alpha^*\beta$?

- Solution 1

PREFIX SEARCH α^* , PREFIX SEARCH $^*\beta$

AND INTERSECT THE RESULTING LIST

EXPENSIVE!

EVEN WORSE IF WE SEARCH FOR SOMETHING

LIKE: $\alpha^*\beta$ AND $\gamma^*\gamma$

LOTS OF COMBINATIONS!!!

- Solution 2

PERMUTEX INDEX

IT TRANSFORMS THE WILDCARD QUERY SO THAT

THE $*$ OCCUR AT THE END.

CONSTRUCTION AND USING THE PERMUTEX INDEX

- ① For every term in our lexicon, index it under all possible rotation of the word with the special char \$

EXAMPLE

Hello is indexed by:

Hello\$, ello\$H, llo\$He, lo\$Hel, o\$Hell, \$Hello

NOTE: Permutex index size $\approx 4 \cdot \text{lexicon size}$

② When querying for $\alpha^* \beta$, the search is reduced to a prefix search of the rotated query $\beta \$ \alpha^*$

How to do it?

- $X \rightarrow \underline{X \$}$ CASE BASE
- $X^* \rightarrow X^* \$ \rightarrow * \$ X \rightarrow \underline{\$ X^*}$
- $*X \rightarrow *X \$ \rightarrow \underline{X \*
- $*X^* \rightarrow *X^* \$ \rightarrow X^* \$^* \rightarrow \underline{X^*}$?
- $X^*Y \rightarrow X^*Y \$ \rightarrow *Y \$ X \rightarrow \underline{Y \$ X^*}$
- $X^*Y^*Z \rightarrow$ lookup for $Y \$ X^*$,
then match with $Z \*