# Ανάκληση Πληροφορίας

# Information Retrieval

Διδάσκων –
Δημήτριος Κατσαρός

Διάλεξη 5η: 26/02/2014

# Phrase queries

# Phrase queries

- Want to answer queries such as "**stanford university**" – as a phrase

- Thus the sentence *"I went to university at Stanford"* is not a match.
  - The concept of phrase queries has proven easily understood by users; about 10% of web queries are phrase queries

- No longer suffices to store only

  *<term : docs>* entries

# A first attempt: Biword indexes

- Index every consecutive pair of terms in the text as a phrase
- For example the text "Friends, Romans, Countrymen" would generate the biwords
  - *friends romans*
  - *romans countrymen*
- Each of these biwords is now a dictionary term
- Two-word phrase query-processing is now immediate.

# Longer phrase queries

- Longer phrases are processed as we did with wild-cards:

- *stanford university palo alto* can be broken into the Boolean query on biwords:

*stanford university AND university palo AND palo alto*

Without the docs, we cannot verify that the docs matching the above Boolean query do contain the phrase.

Can have false positives!

# Extended biwords

- Parse the indexed text and perform part-of-speech-tagging (POST).
- Bucket the terms into (say) Nouns (N) and articles/prepositions (X).
- Now deem any string of terms of the form NX*N to be an <u>extended biword</u>.
  - Each such extended biword is now made a term in the dictionary.
- Example:  ***catcher in the rye***
  - **N        X  X   N**
- Query processing: parse it into N's and X's
  - Segment query into enhanced biwords
  - Look up index

# Issues for biword indexes

- False positives, as noted before
- Index blowup due to bigger dictionary

- For extended biword index, parsing longer queries into conjunctions:
  - E.g., the query **tangerine trees and marmalade skies** is parsed into
  - **tangerine trees** AND **trees and marmalade** AND **marmalade skies**

- Not standard solution (for all biwords)

# Solution 2: Positional indexes

- Store, for each **term**, entries of the form:

  <number of docs containing **term**;

  *doc1*: position1, position2 … ;

  *doc2*: position1, position2 … ;
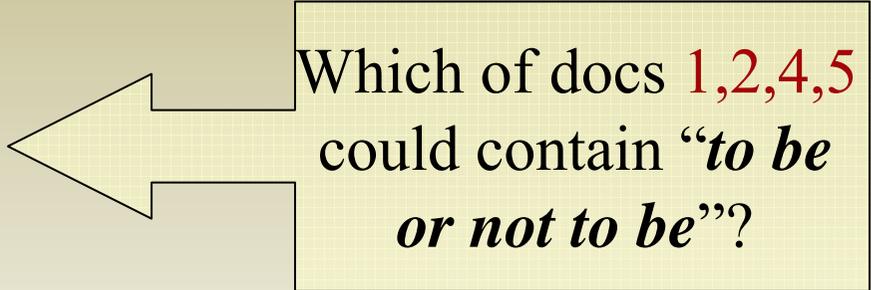
  etc.>

# Positional index example

<be: 993427;
1: 7, 18, 33, 72, 86, 231;
2: 3, 149;
4: 17, 191, 291, 430, 434;
5: 363, 367, …>

Which of docs 1,2,4,5 could contain "*to be or not to be*"?

- Can compress position values/offsets
- Nevertheless, this expands postings storage *substantially*

# Processing a phrase query

- Extract inverted index entries for each distinct term: **to, be, or, not.**
- Merge their *doc:position* lists to enumerate all positions with "**to be or not to be**".

    - *to:*
        - *2*:1,17,74,222,551; *4*:8,16,190,429,433; *7*:13,23,191; ...
    - *be:*
        - *1*:17,19; *4*:17,191,291,430,434; *5*:14,19,101; ...

- Same general method for proximity searches

# Proximity queries

- LIMIT! /3 STATUTE /3 FEDERAL /2 TORT Here, /$k$ means "within $k$ words of".

- Clearly, positional indexes can be used for such queries; biword indexes cannot.

- Exercise: Adapt the linear merge of postings to handle proximity queries. Can you make it work for any value of $k$?

# Positional index size

- You can compress position values/offsets: we'll talk bout that in lecture 5

- Nevertheless, a positional index expands postings storage *substantially*

- Nevertheless, it is now standardly used because of the power and usefulness of phrase and proximity queries … whether used explicitly or implicitly in a ranking retrieval system.

# Positional index size

- Need an entry for each occurrence, not just once per document
- Index size depends on average document size
  - Average web page has <1000 terms
  - SEC filings, books, even some epic poems … easily 100,000 terms
- Consider a term with frequency 0.1%

**Why?**

| Document size | Postings | Positional postings |
|---|---|---|
| 1000 | 1 | 1 |
| 100,000 | 1 | 100 |

# Rules of thumb

- A positional index is 2–4 as large as a non-positional index

- Positional index size 35–50% of volume of original text

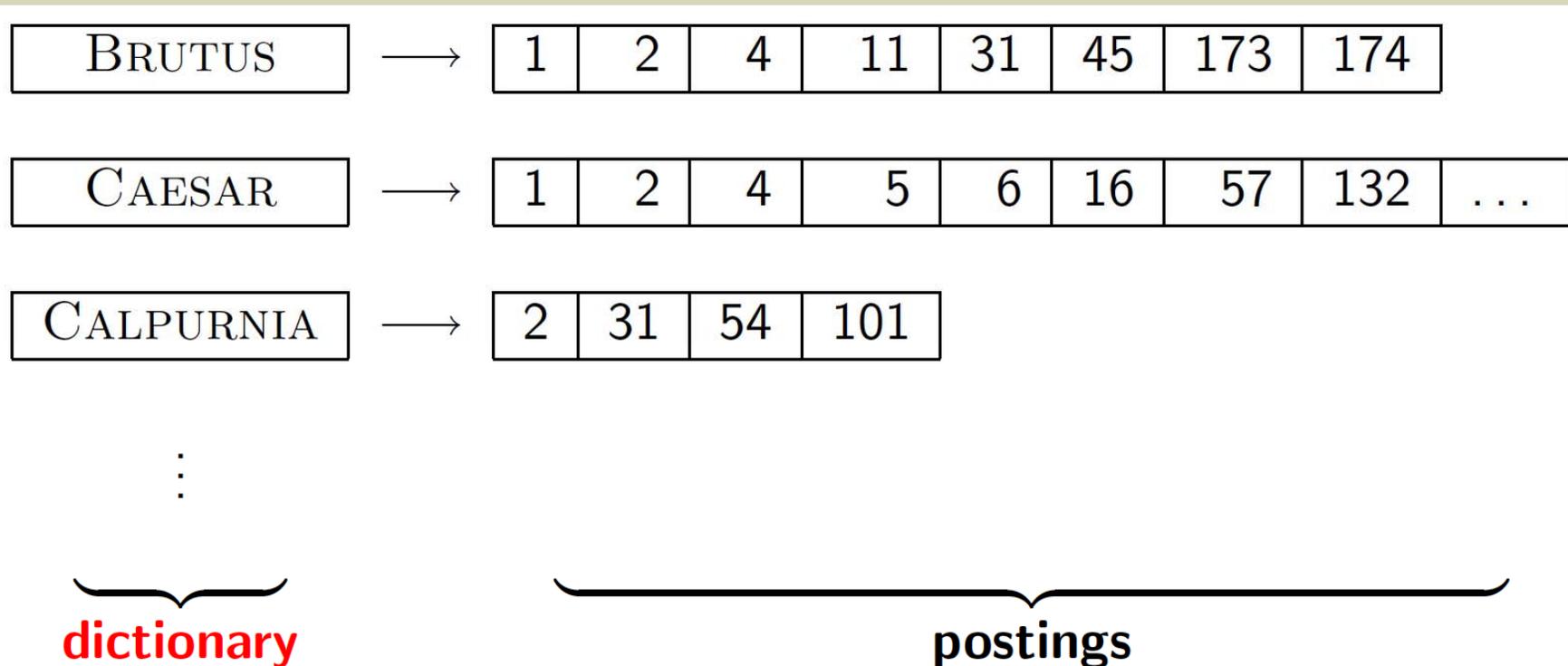- Caveat: all of this holds for "English-like" languages

# Combination schemes

- These two approaches can be profitably combined
  - For particular phrases (**"Michael Jackson", "Britney Spears"**) it is inefficient to keep on merging positional postings lists
    - Even more so for phrases like **"The Who"**
- Williams et al. (2004) evaluate a more sophisticated mixed indexing scheme
  - A typical web query mixture was executed in ¼ of the time of using just a positional index
  - It required 26% more space than having a positional index alone

# Dictionary data structures for inverted indexes

- The dictionary data structure stores the term vocabulary, document frequency, pointers to each postings list … in what data structure?

| BRUTUS | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|--------|---|---|---|---|----|----|----|-----|-----|

| CAESAR | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | … |
|--------|---|---|---|---|---|---|----|----|-----|---|

| CALPURNIA | → | 2 | 31 | 54 | 101 |
|-----------|---|---|----|----|-----|

**dictionary**          **postings**

# A naïve dictionary

- An array of struct:

| term | document frequency | pointer to postings list |
|------|--------------------|--------------------------|
| a | 656,265 | ⟶ |
| aachen | 65 | ⟶ |
| ... | ... | ... |
| zulu | 221 | ⟶ |

char[20]   int                Postings *
20 bytes   4/8 bytes          4/8 bytes

- How do we store a dictionary in memory efficiently?
- How do we quickly look up elements at query time?

# Dictionary data structures

- Two main choices:
  - Hash table
  - Tree
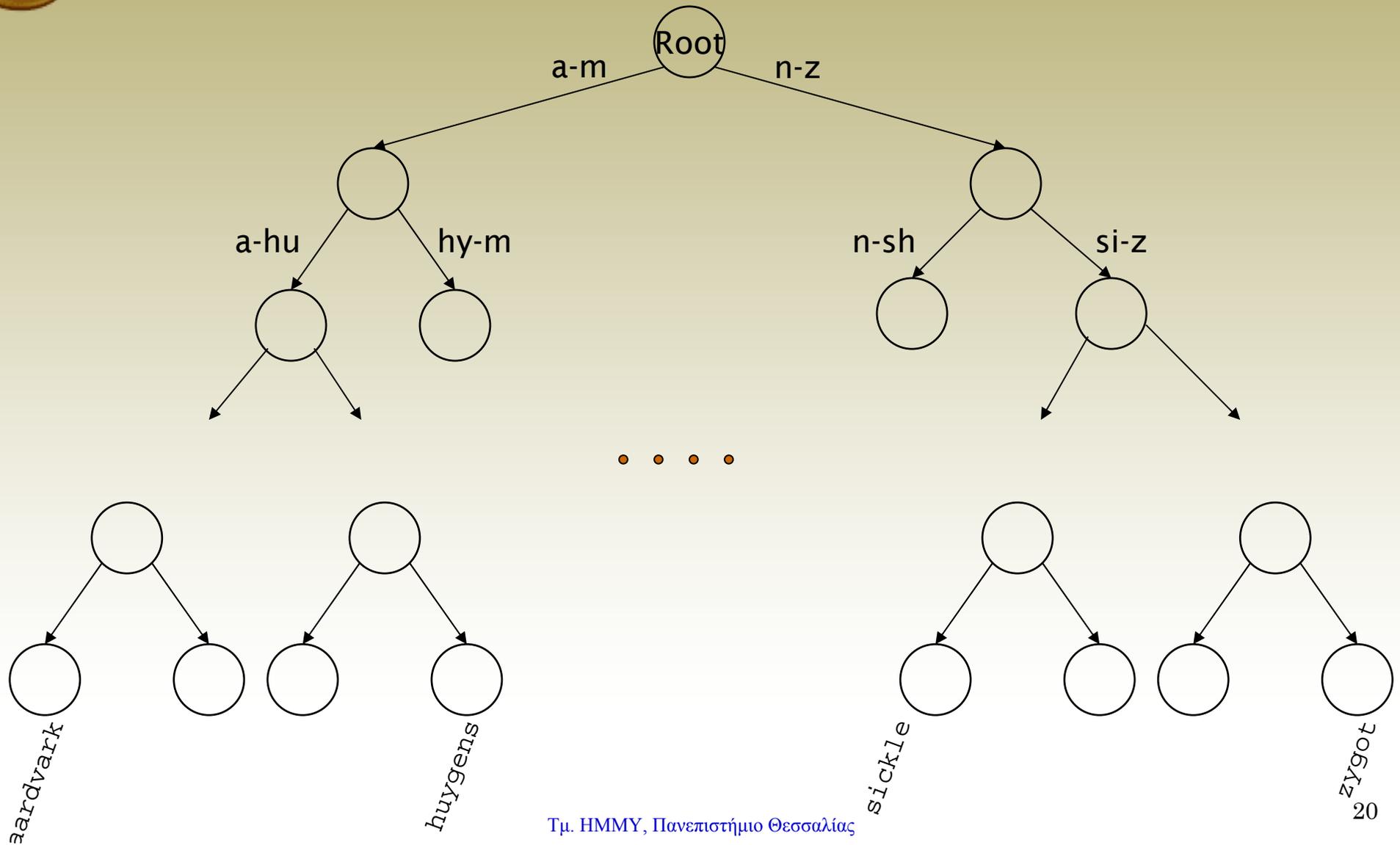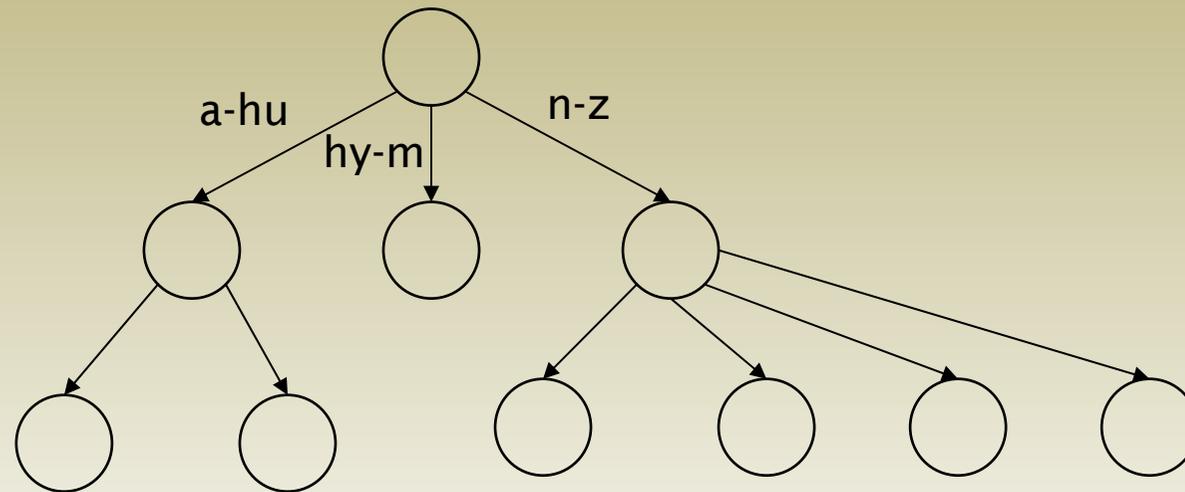- Some IR systems use hashes, some trees

# Hashes

- Each vocabulary term is hashed to an integer
  - (We assume you've seen hashtables before)
- Pros:
  - Lookup is faster than for a tree: O(1)
- Cons:
  - No easy way to find minor variants:
    - judgment/judgement
  - No prefix search          [tolerant retrieval]
  - If vocabulary keeps growing, need to occasionally do the expensive operation of rehashing *everything*

# Tree: binary tree



Root

a-m    n-z

a-hu    hy-m    n-sh    si-z

• • • •

aardvark    huygens    sickle    zygot

# Tree: B-tree



- Definition: Every internal nodel has a number of children in the interval [a,b] where *a, b* are appropriate natural numbers, e.g., [2,4].

# Trees

- Simplest: binary tree
- More usual: B-trees
- Trees require a standard ordering of characters and hence strings … but we standardly have one
- Pros:
  - Solves the prefix problem (terms starting with *hyp*)
- Cons:
  - Slower: $O(\log M)$  [and this requires *balanced* tree]
  - Rebalancing binary trees is expensive
    - But B-trees mitigate the rebalancing problem