# CS60092: Information Retrieval
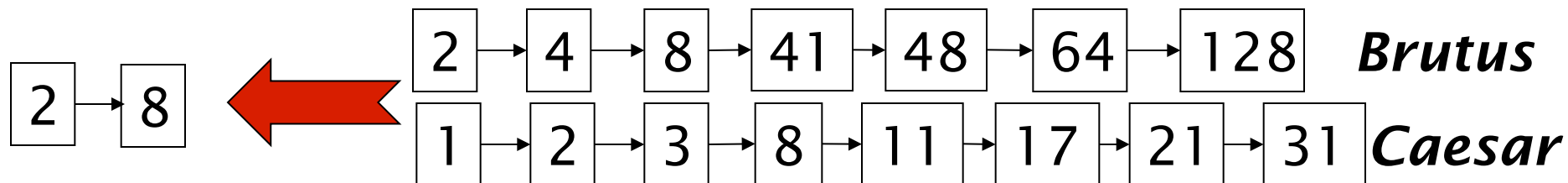
Sourangshu Bhattacharya

# Introduction to
# **Information Retrieval**

Faster postings merges:
Skip pointers/Skip lists

# Recall basic merge

- Walk through the two postings simultaneously, in time linear in the total number of postings entries
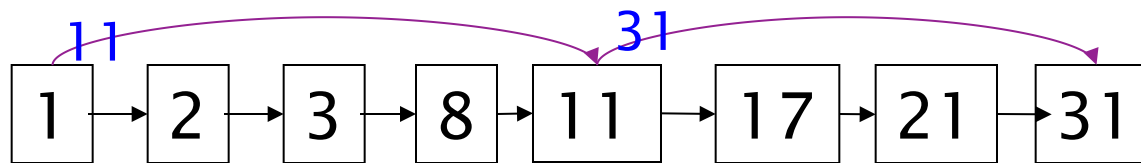
| 2 | 4 | 8 | 41 | 48 | 64 | 128 | ***Brutus*** |

| 2 | 8 | ← |

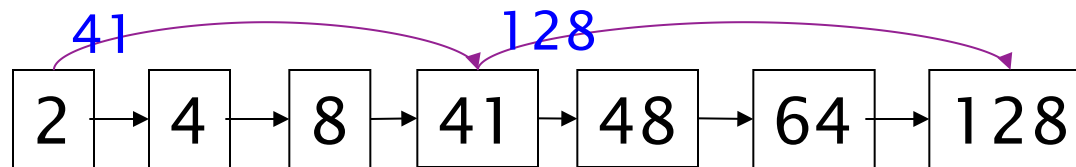| 1 | 2 | 3 | 8 | 11 | 17 | 21 | 31 | ***Caesar*** |

If the list lengths are *m* and *n*, the merge takes O(*m+n*) operations.
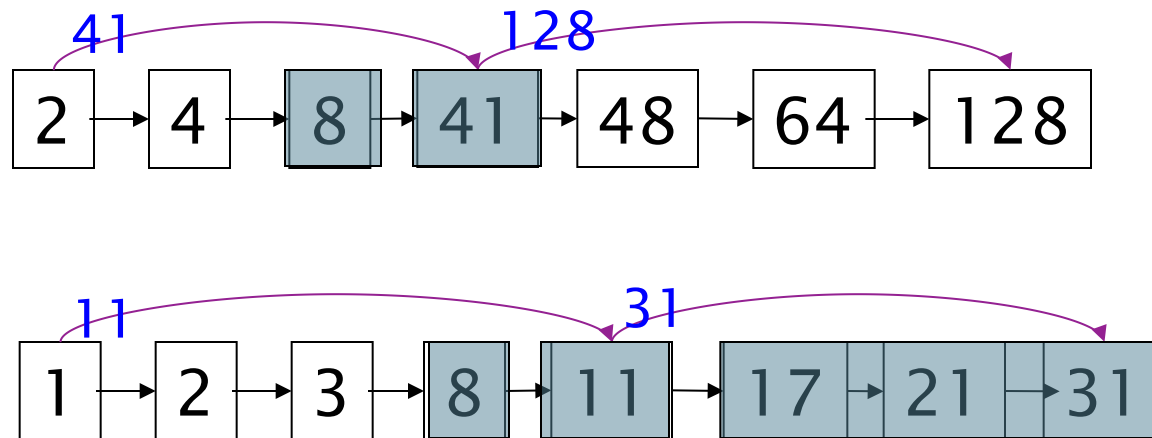
Can we do better?
Yes (if the index isn't changing too fast).

# Augment postings with skip pointers (at indexing time)



- Why?
- <u>To skip postings that will not figure in the search results.</u>
- How?
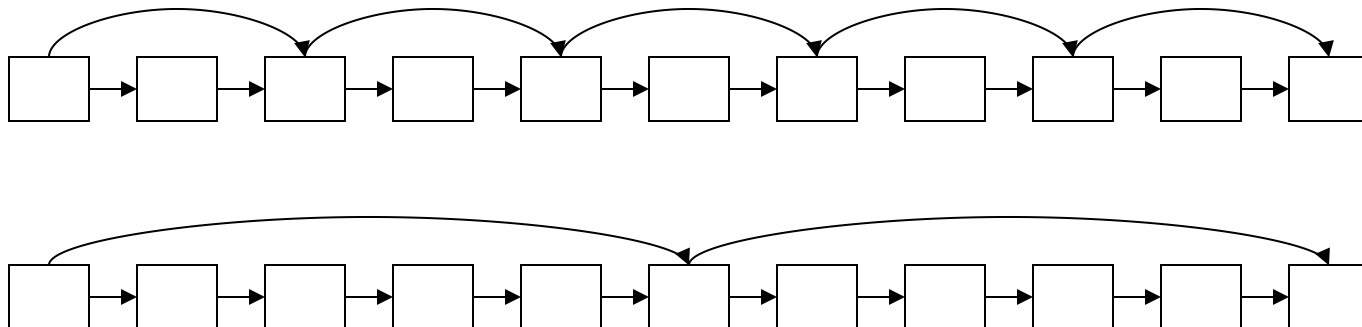- Where do we place skip pointers?

# Query processing with skip pointers



Suppose we've stepped through the lists until we process **8** on each list. We match it and advance.

We then have **41** and **11** on the lower. **11** is smaller.

But the skip successor of **11** on the lower list is **31**, so we can skip ahead past the intervening postings.

# Where do we place skips?

- Tradeoff:
    - More skips → shorter skip spans ⇒ more likely to skip. But lots of comparisons to skip pointers.
    - Fewer skips → few pointer comparison, but then long skip spans ⇒ few successful skips

# Placing skips

- Simple heuristic: for postings of length $L$, use $\sqrt{L}$ evenly-spaced skip pointers   [Moffat and Zobel 1996]

- Easy if the index is relatively static; harder if $L$ keeps changing because of updates.

- This definitely used to help; with modern hardware it may not unless you're memory-based   [Bahle et al. 2002]

  - The I/O cost of loading a bigger postings list can outweigh the gains from quicker in memory merging!

# Phrase queries

- We want to answer a query such as [stanford university] – as a phrase.

- Thus *The inventor Stanford Ovshinsky never went to university* should not be a match.

- The concept of phrase query has proven easily understood by users.

- About 10% of web queries are phrase queries.

- Consequence for inverted index: it no longer suffices to store docIDs in postings lists for terms.

- Two ways of extending the inverted index:
  - biword index
  - positional index

# Biword indexes

▪Index every consecutive pair of terms in the text as a phrase.

▪For example, *Friends, Romans, Countrymen* would generate two biwords: *"friends romans"* and *"romans countrymen"*

▪Each of these biwords is now a vocabulary term.

▪Two-word phrases can now easily be answered.

# Longer phrase queries

▪A long phrase like *"stanford university palo alto"* can be represented as the Boolean query "STANFORD UNIVERSITY" AND "UNIVERSITY PALO" AND "PALO ALTO"

▪Does this always guarantee the correct match? -- We need to do post-filtering of hits to identify subset that actually contains the 4-word phrase.

▪What about phrases like, *"abolition of slavery"*?

# Issues with biword indexes

- Why are biword indexes rarely used?
- False positives, as noted above
- Index blowup due to very large term vocabulary

- *What can be an alternative?*

# Positional indexes

- Positional indexes are a more efficient alternative to biword indexes.

- Postings lists in a nonpositional index: each posting is just a docID

- Postings lists in a positional index: each posting is a docID and a list of positions

# Positional indexes: Example

Query: *"to$_1$ be$_2$ or$_3$ not$_4$ to$_5$ be$_6$"* TO, 993427:

‹ 1: ‹7, 18, 33, 72, 86, 231›;

2: ‹1, 17, 74, 222, 255›;

4: ‹8, 16, 190, 429, 433›;

5: ‹363, 367›;

7: ‹13, 23, 191›; . . . ›

BE, 178239:

‹ 1: ‹17, 25›;

4: ‹17, 191, 291, 430, 434›;

5: ‹14, 19, 101›; . . . › Document 4 is a match!

# Proximity search

- We just saw how to use a positional index for phrase searches.

- *Can we also use it for proximity search?*

- For example: employment /4 place

- Find all documents that contain EMPLOYMENT and PLACE within 4 words of each other.

- *Employment agencies that place healthcare workers are seeing growth* is a hit.

- *Employment agencies that have learned to adapt now place healthcare workers* is not a hit.

# Proximity search

- Use the positional index

- Simplest algorithm: look at cross-product of positions of (i) EMPLOYMENT in document and (ii) PLACE in document

- Very inefficient for frequent words, especially stop words

- Note that we want to return the actual matching positions, not just a list of documents.

# Combination scheme

- Biword indexes and positional indexes can be profitably combined.

- Many biwords are extremely frequent: Michael Jackson etc

- For these biwords, increased speed compared to positional postings intersection is substantial.

- Combination scheme: Include frequent biwords as vocabulary terms in the index. Do all other phrases by positional intersection.

- Williams et al. (2004) evaluate a more sophisticated mixed indexing scheme – *Next Word Index*. Faster than a positional index, at a cost of 26% more space for index.

# Introduction to
# **Information Retrieval**

## Dictionaries, Tolerant Retrieval

# 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 | $\longrightarrow$ |
| aachen | 65 | $\longrightarrow$ |
| . . . | . . . | . . . |
| zulu | 221 | $\longrightarrow$ |

char[20]   int            Postings *

20 bytes   4/8 bytes      4/8 bytes

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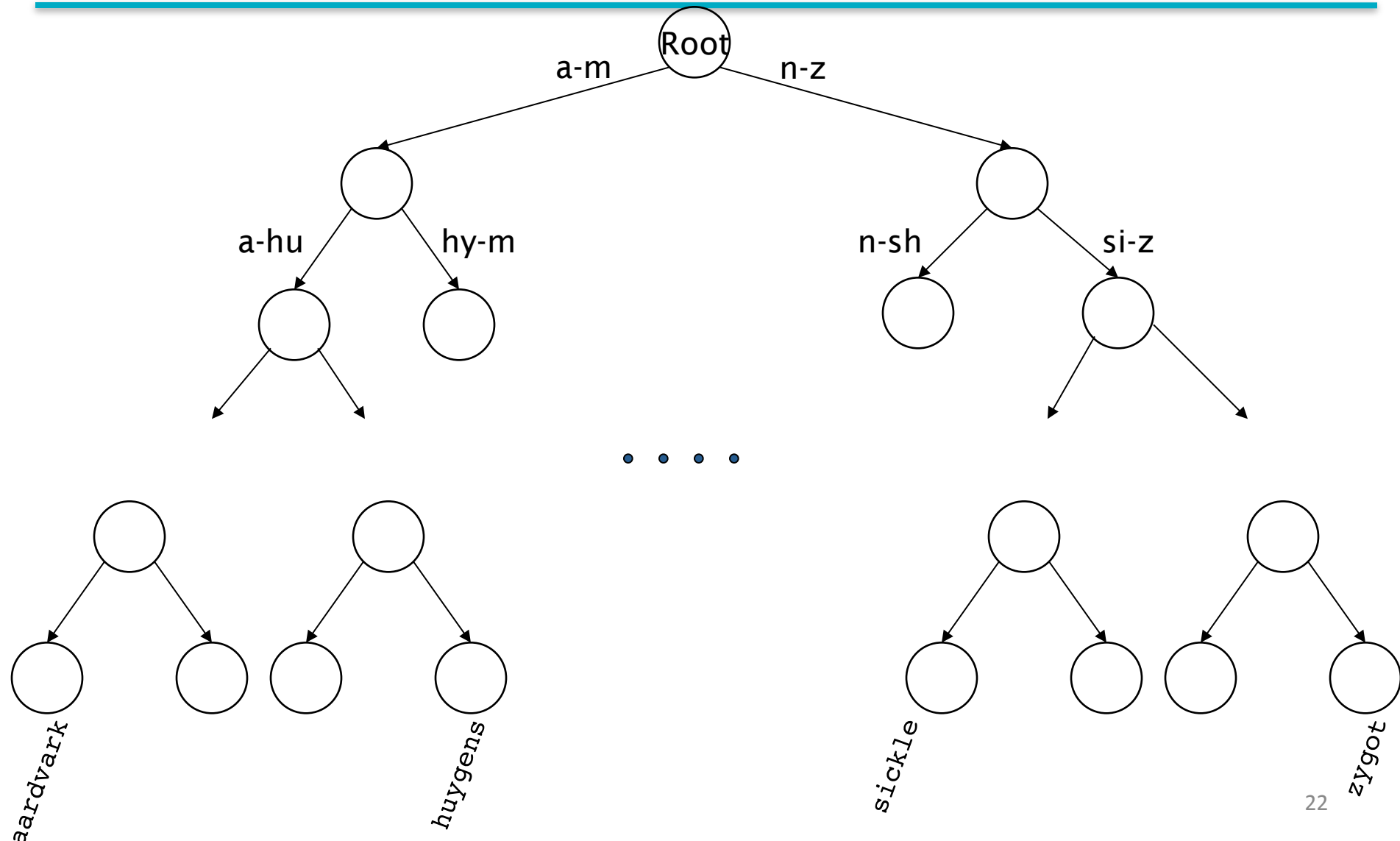- How do we quickly look up elements at query time?

# Dictionary data structures

- Two main choices:
    - Hashtables
    - Trees
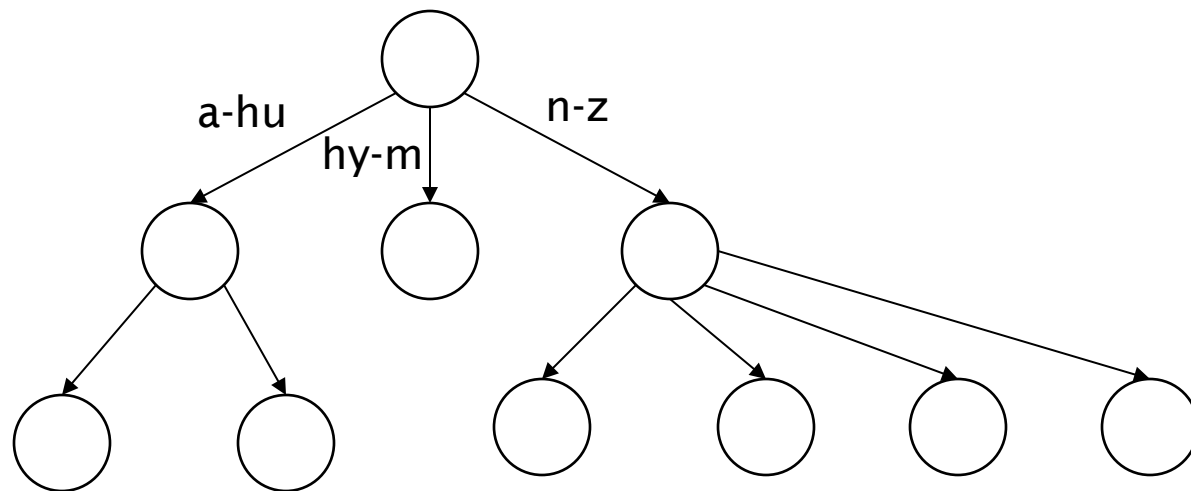
- Some IR systems use hashtables, some trees

# Hashtables

- 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

# 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 typically 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

# WILD-CARD QUERIES

# Wild-card queries: *

- ***mon*:** find all docs containing any word beginning with "mon".

- Easy with binary tree (or B-tree) lexicon: retrieve all words in range: ***mon ≤ w < moo***

- ***mon:** find words ending in "mon": harder

  - Maintain an additional B-tree for terms *backwards.*

  Can retrieve all words in range: ***nom ≤ w < non**.*

Exercise: from this, how can we enumerate all terms meeting the wild-card query ***pro*cent*** ?

# Query processing

- At this point, we have an enumeration of all terms in the dictionary that match the wild-card query.

- We still have to look up the postings for each enumerated term.

- E.g., consider the query:

  ***se*ate** *AND* **fil*er***

  This may result in the execution of many Boolean *AND* queries.

# B-trees handle *'s at the end of a query term

- How can we handle *'s in the middle of query term?
  - ***co\*tion***
- We could look up ***co\**** AND ***\*tion*** in a B-tree and intersect the two term sets
  - Expensive
- The solution: transform wild-card queries so that the *'s occur at the end
- This gives rise to the **Permuterm** Index.

# Permuterm index

- For term **_hello_**, index under:
  - **_hello$, ello$h, llo$he, lo$hel, o$hell, $hello_**

  where $ is a special symbol.

- Queries:
  - **X**   lookup on **X$**         **X\***   lookup on   $**X\***
  - **\*X**   lookup on **X$\***      **\*X\***  lookup on   **X\***
  - **X\*Y** lookup on **Y$X\***         **X\*Y\*Z**     ??? Exercise!

Query = **_hel\*o_**
**X=_hel_, **Y=_o_**
Lookup **_o$hel\*_**

29

# Permuterm query processing

- Rotate query wild-card to the right

- Now use B-tree lookup as before.

- *Permuterm problem: ≈ quadruples lexicon size*

Empirical observation for English.

# Bigram (*k*-gram) indexes

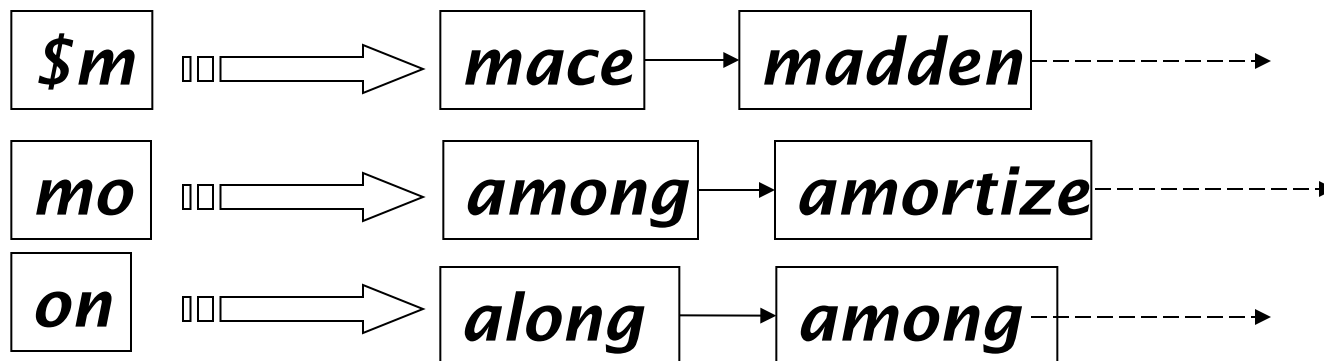- Enumerate all *k*-grams (sequence of *k* chars) occurring in any term

- *e.g.,* from text "***April is the cruelest month***" we get the 2-grams (*bigrams*)

  $a,ap,pr,ri,il,l$,$i,is,s$,$t,th,he,e$,$c,cr,ru,
  ue,el,le,es,st,t$, $m,mo,on,nt,h$

  - $ is a special word boundary symbol

- Maintain a <u>*second*</u> inverted index <u>*from bigrams to dictionary terms*</u> that match each bigram.

# Bigram index example

- The *k*-gram index finds *terms* based on a query consisting of *k*-grams (here *k*=2).

| $m → mace → madden ------→

mo → among → amortize ------→

on → along → among ------→

# Processing wild-cards

- Query ***mon\**** can now be run as
  - ***$m*** *AND* ***mo*** *AND* ***on***
- Gets terms that match AND version of our wildcard query.
- But we'd enumerate ***moon***.
- Must post-filter these terms against query.
- Surviving enumerated terms are then looked up in the term-document inverted index.
- Fast, space efficient (compared to permuterm).

# Processing wild-card queries

- As before, we must execute a Boolean query for each enumerated, filtered term.

- Wild-cards can result in expensive query execution (very large disjunctions…)

  - pyth* AND prog*

- If you encourage "laziness" people will respond!

| |
|---|
| Search |

Type your search terms, use '*' if you need to.
E.g., Alex* will match Alexander.

# Recap

- Dictionary Data Structures: Hashtables, B-trees (why not tries?)

- Tolerant Retrieval

  - Wildcard queries (permuterm index, k-gram index)

  - Spell correction contd ...

# SPELLING CORRECTION

# Spell correction

- Two principal uses
  - Correcting document(s) being indexed
  - Correcting user queries to retrieve "right" answers
- Two main flavors:
  - Isolated word
    - Check each word on its own for misspelling: *jacson*
    - Will not catch typos resulting in correctly spelled words
    - e.g., *from → form*
  - Context-sensitive
    - Look at surrounding words,
    - e.g., *I flew <u>form</u> Heathrow to Narita.*

37

# Document correction

- Especially needed for OCR'ed documents
  - Correction algorithms are tuned for this: rn/m
  - Can use domain-specific knowledge
    - E.g., OCR can confuse O and D more often than it would confuse O and I (adjacent on the QWERTY keyboard, so more likely interchanged in typing).
- But also: web pages and even printed material have typos
- Goal: the dictionary contains fewer misspellings
- But often we don't change the documents and instead fix the query-document mapping

# Query mis-spellings

- Our principal focus here
  - E.g., the query **_Alanis Morisett_**
- We can either
  - Retrieve documents indexed by the correct spelling, OR
  - Return several suggested alternative queries with the correct spelling
    - *Did you mean … ?*

# Isolated word correction

- Fundamental premise – there is a lexicon from which the correct spellings come

- Two basic choices for this
  - A standard lexicon such as
    - Webster's English Dictionary
    - An "industry-specific" lexicon – hand-maintained
  - The lexicon of the indexed corpus
    - E.g., all words on the web
    - All names, acronyms etc.
    - (Including the mis-spellings)

# Isolated word correction

- Given a lexicon and a character sequence Q, return the words in the lexicon closest to Q

- What's "closest"?

- We'll study several alternatives
  - Edit distance (Levenshtein distance)
  - Weighted edit distance
  - *n*-gram overlap

# Edit distance

- Given two strings $S_1$ and $S_2$, the minimum number of operations to convert one to the other

- Operations are typically character-level

  - Insert, Delete, Replace, (Transposition)

- E.g., the edit distance from **dof** to **dog** is 1

  - From **cat** to **act** is 2    (Just 1 with transpose.)

  - from **cat** to **dog** is 3.

- Generally found by dynamic programming.

# Edit distance

- Minimum number of edit operations needed to change string $s_1$ to string $s_2$.

- When the edit operations are:

  - Add

  - Delete

  - Replace

- The edit distance is called Levenstein distance.

# Edit distance

$\text{EDITDISTANCE}(s_1, s_2)$
1   $int\ m[i,j] = 0$
2   **for** $i \leftarrow 1$ **to** $|s_1|$
3   **do** $m[i, 0] = i$
4   **for** $j \leftarrow 1$ **to** $|s_2|$
5   **do** $m[0, j] = j$
6   **for** $i \leftarrow 1$ **to** $|s_1|$
7   **do for** $j \leftarrow 1$ **to** $|s_2|$
8      **do** $m[i, j] = \min\{m[i-1, j-1] + \text{if } (s_1[i] = s_2[j]) \text{ then } 0 \text{ else } 1\text{fi},$
9                          $m[i-1, j] + 1,$
10                         $m[i, j-1] + 1\}$
11  **return** $m[|s_1|, |s_2|]$

# Edit Distance

|   |   |   | f |   | a |   | s |   | t |   |
|---|---|---|---|---|---|---|---|---|---|---|
|   |   | 0 | 1 | 1 | 2 | 2 | 3 | 3 | 4 | 4 |
| c |   | **1** / **1** | *1* / *2* | 2 / *1* | **2** / **2** | 3 / **2** | **3** / **3** | 4 / **3** | **4** / **4** | 5 / **4** |
| a |   | **2** / **2** | **2** / *3* | 2 / **2** | *1* / *3* | 3 / *1* | *3* / *2* | 4 / *2* | 4 / *3* | 5 / *3* |
| t |   | **3** / **3** | **3** / *4* | 3 / **3** | 3 / *4* | **2** / **2** | **2** / *3* | 3 / *2* | **2** / *3* | 4 / **2** |
| s |   | **4** / **4** | **4** / *5* | 4 / **4** | 4 / *5* | **3** / **3** | **2** / *4* | 3 / **2** | *3* / *3* | *3* / *3* |

# Weighted edit distance

- As above, but the weight of an operation depends on the character(s) involved

  - Meant to capture OCR or keyboard errors
    Example: *m* more likely to be mis-typed as *n* than as *q*

  - Therefore, replacing *m* by *n* is a smaller edit distance than by *q*

  - This may be formulated as a probability model

- Requires weight matrix as input

- Modify dynamic programming to handle weights

# Using edit distances for correction

- Given query, first enumerate all character sequences within a preset (weighted) edit distance (e.g., 2)

- Intersect this set with list of "correct" words

- Show terms you found to user as suggestions

- Alternatively,

  - We can look up all possible corrections in our inverted index and return all docs … slow

  - We can run with a single most likely correction

- The alternatives disempower the user, but save a round of interaction with the user

# *n*-gram overlap

- Enumerate all the *n*-grams in the query string as well as in the lexicon

- Use the *n*-gram index (recall wild-card search) to retrieve all lexicon terms matching any of the query *n*-grams

- Threshold by number of matching *n*-grams
  - Variants – weight by keyboard layout, etc.

# Example with trigrams

- Suppose the text is ***november***
  - Trigrams are *nov, ove, vem, emb, mbe, ber*.
- The query is ***december***
  - Trigrams are *dec, ece, cem, emb, mbe, ber*.
- So 3 trigrams overlap (of 6 in each term)
- How can we turn this into a normalized measure of overlap?
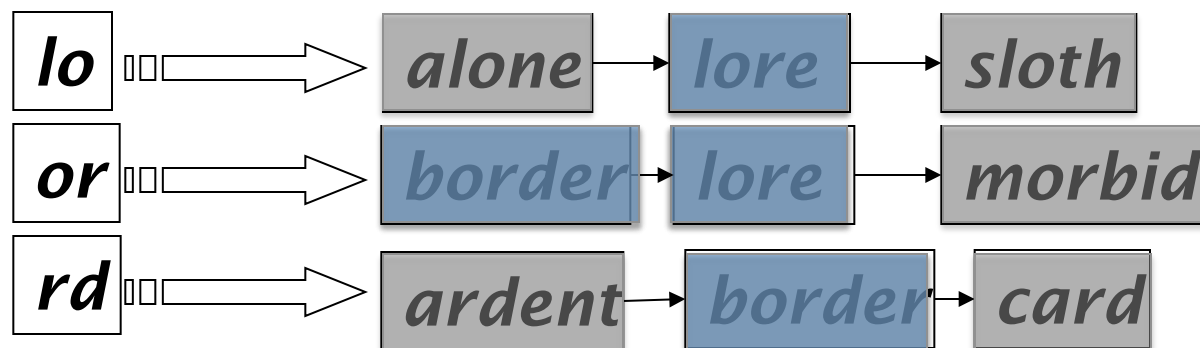
# One option – Jaccard coefficient

- A commonly-used measure of overlap
- Let *X* and *Y* be two sets; then the J.C. is

$$|X \cap Y| / |X \cup Y|$$

- Equals 1 when *X* and *Y* have the same elements and zero when they are disjoint
- *X* and *Y* don't have to be of the same size
- Always assigns a number between 0 and 1
  - Now threshold to decide if you have a match
  - E.g., if J.C. > 0.8, declare a match

# Matching trigrams

- Consider the query *lord* – we wish to identify words matching 2 of its 3 bigrams (*lo, or, rd*)



| *lo* | → | *alone* → *lore* → *sloth* |
| *or* | → | *border* → *lore* → *morbid* |
| *rd* | → | *ardent* → *border* → *card* |

Standard postings "merge" will enumerate …

Adapt this to using Jaccard (or another) measure.

# Context-sensitive spell correction

- Text: ***I flew <u>from </u>Heathrow to Narita.***

- Consider the phrase query ***"flew <u>form </u>Heathrow"***

- We'd like to respond

   Did you mean "***flew from Heathrow***"?

because no docs matched the query phrase.

# Context-sensitive correction

- Need surrounding context to catch this.

- First idea: retrieve dictionary terms close (in weighted edit distance) to each query term

- Now try all possible resulting phrases with one word "fixed" at a time
  - *flew from heathrow*
  - *fled form heathrow*
  - *flea form heathrow*

- **Hit-based spelling correction:** Suggest the alternative that has lots of hits.

# Exercise

- Suppose that for ***"flew <u>form </u>Heathrow"*** we have 7 alternatives for flew, 19 for form and 3 for heathrow.

How many "corrected" phrases will we enumerate in this scheme?

# Another approach

- Break phrase query into a conjunction of biwords (Lecture 2).

- Look for biwords that need only one term corrected.

- Enumerate only phrases containing "common" biwords.

# SOUNDEX

# Soundex

- Class of heuristics to expand a query into phonetic equivalents

    - Language specific – mainly for names :

    - E.g., ***chebyshev*** → ***tchebycheff***

- Invented for the U.S. census … in 1918

# Soundex – typical algorithm

- Turn every token to be indexed into a 4-character reduced form

- Do the same with query terms

- Build and search an index on the reduced forms
  - (when the query calls for a soundex match)

- http://www.creativyst.com/Doc/Articles/SoundEx1/SoundEx1.htm#Top

# Soundex – typical algorithm

1. Retain the first letter of the word.
2. Change all occurrences of the following letters to '0' (zero):
   'A', E', 'I', 'O', 'U', 'H', 'W', 'Y'.
3. Change letters to digits as follows:

- B, F, P, V → 1
- C, G, J, K, Q, S, X, Z → 2
- D,T → 3
- L → 4
- M, N → 5
- R → 6

# Soundex continued

4.  Remove all pairs of consecutive digits.

5.  Remove all zeros from the resulting string.

6.  Pad the resulting string with trailing zeros and return the first four positions, which will be of the form <uppercase letter> <digit> <digit> <digit>.

E.g., **Herman** becomes H655.

Will **hermann** generate the same code?

# Soundex

- Soundex is the classic algorithm, provided by most databases (Oracle, Microsoft, …)

- How useful is soundex?

- Not very – for information retrieval

- Okay for "high recall" tasks (e.g., Interpol), though biased to names of certain nationalities

- Zobel and Dart (1996) show that other algorithms for phonetic matching perform much better in the context of IR

# What queries can we process?

- ■ We have
  - ■ Positional inverted index with skip pointers
  - ■ Wild-card index
  - ■ Spell-correction
  - ■ Soundex
- ■ Queries such as

**(SPELL(moriset) /3 toron\*to) OR** *SOUNDEX***(chaikofski)**