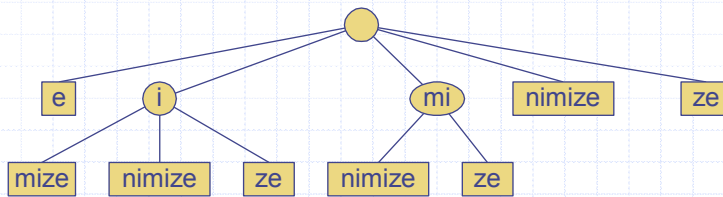


Presentation for use with the textbook [Data Structures and Algorithms in Java, 6th edition](#), by M. T. Goodrich, R. Tamassia, and M. H. Goldwasser, Wiley, 2014

Tries

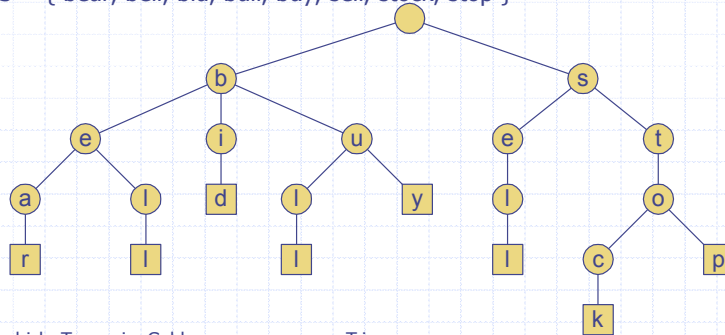


Preprocessing Strings

- ◆ Preprocessing the pattern speeds up pattern matching queries
 - After preprocessing the pattern, KMP's algorithm performs pattern matching in time proportional to the text size
- ◆ If the text is large, immutable and searched for often (e.g., works by Shakespeare), we may want to preprocess the text instead of the pattern
- ◆ A trie is a compact data structure for representing a set of strings, such as all the words in a text
 - A trie supports pattern matching queries in time proportional to the pattern size

Standard Tries

- ◆ The standard trie for a set of strings S is an ordered tree such that:
 - Each node but the root is labeled with a character
 - The children of a node are alphabetically ordered
 - The paths from the external nodes to the root yield the strings of S
- ◆ Example: standard trie for the set of strings $S = \{ \text{bear, bell, bid, bull, buy, sell, stock, stop} \}$



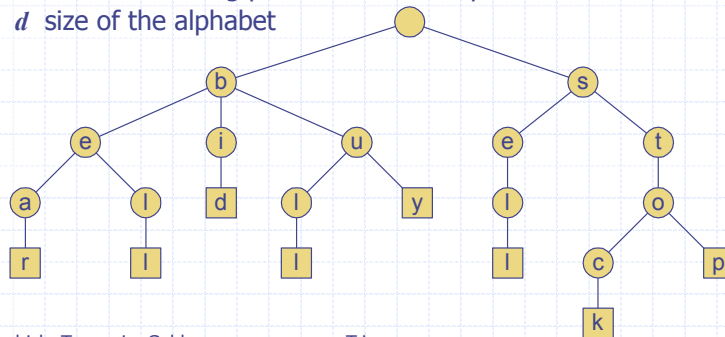
© 2014 Goodrich, Tamassia, Goldwasser

Tries

3

Analysis of Standard Tries

- ◆ A standard trie uses $O(n)$ space and supports searches, insertions and deletions in time $O(dm)$, where:
 - n total size of the strings in S
 - m size of the string parameter of the operation
 - d size of the alphabet



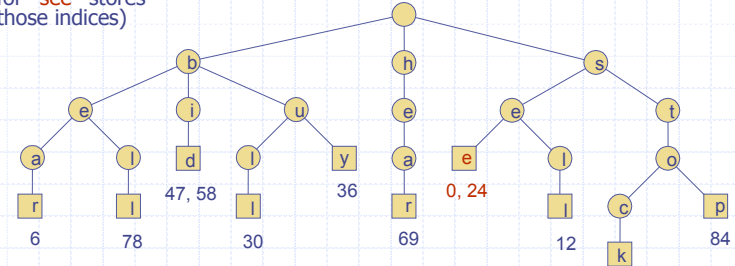
© 2014 Goodrich, Tamassia, Goldwasser

Tries

4

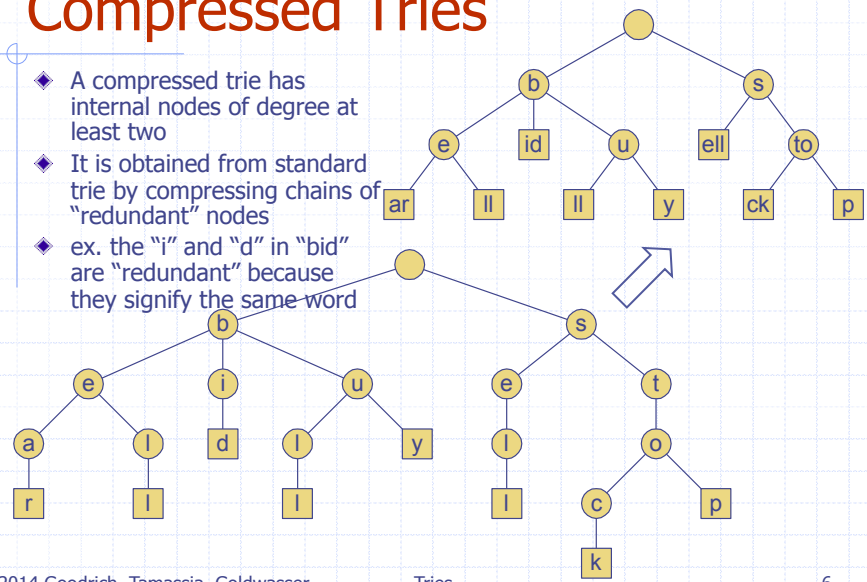
Word Matching with a Trie

- ◆ insert the words of the text into trie
- ◆ Each leaf is associated w/ one particular word
- ◆ leaf stores indices where associated word begins ("see" starts at index 0 & 24, leaf for "see" stores those indices)



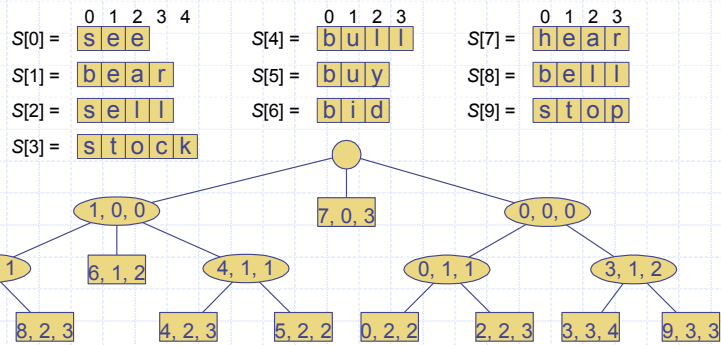
Compressed Tries

- ◆ A compressed trie has internal nodes of degree at least two
- ◆ It is obtained from standard trie by compressing chains of "redundant" nodes
- ◆ ex. the "i" and "d" in "bid" are "redundant" because they signify the same word



Compact Representation

- ◆ Compact representation of a compressed trie for an array of strings:
 - Stores at the nodes ranges of indices instead of substrings
 - Uses $O(s)$ space, where s is the number of strings in the array
 - Serves as an auxiliary index structure



© 2014 Goodrich, Tamassia, Goldwasser

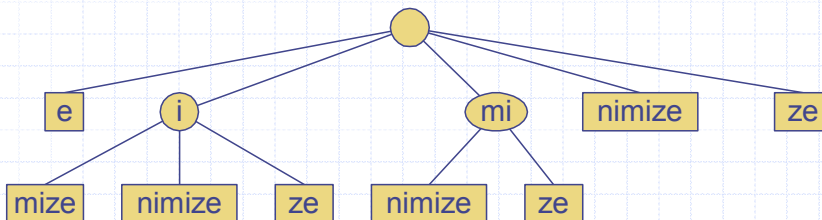
Tries

7

Suffix Trie

- ◆ The suffix trie of a string X is the compressed trie of all the suffixes of X

m i n i m i z e
0 1 2 3 4 5 6 7



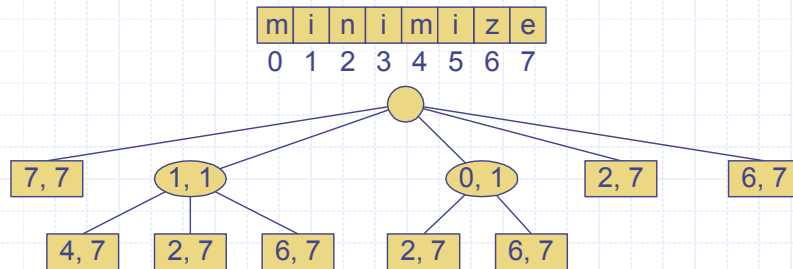
© 2014 Goodrich, Tamassia, Goldwasser

Tries

8

Analysis of Suffix Tries

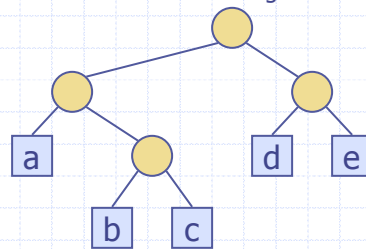
- ◆ Compact representation of the suffix trie for a string X of size n from an alphabet of size d
 - Uses $O(n)$ space
 - Supports arbitrary pattern matching queries in X in $O(dm)$ time, where m is the size of the pattern
 - Can be constructed in $O(n)$ time



Encoding Trie (1)

- ◆ A code is a mapping of each character of an alphabet to a binary code-word
- ◆ A prefix code is a binary code such that no code-word is the prefix of another code-word
- ◆ An encoding trie represents a prefix code
 - Each leaf stores a character
 - The code word of a character is given by the path from the root to the leaf storing the character (0 for a left child and 1 for a right child)

00	010	011	10	11
a	b	c	d	e



Encoding Trie (2)

- ◆ Given a text string X , we want to find a prefix code for the characters of X that yields a small encoding for X
 - Frequent characters should have short code-words
 - Rare characters should have long code-words
- ◆ Example
 - $X = \text{abracadabra}$
 - T_1 encodes X into 29 bits
 - T_2 encodes X into 24 bits

