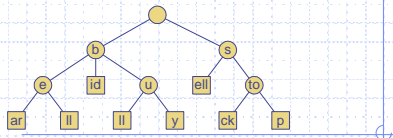


Tries



Preprocessing Strings

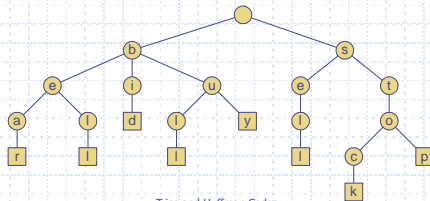
- ◆ A trie (retrieval) is a special kind information access tree.
- ◆ If the text is large, immutable and searched for often (e.g., works by Shakespeare), we may want to preprocess the text.
- ◆ A trie is a compact data structure for representing a set of strings, such as all the words in a text
 - A trie supports pattern matching queries in time proportional to the pattern size

Tries and Huffman Codes

2

Standard Trie

- ◆ 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 root to external nodes yield the strings of S
- ◆ Example: standard trie for the set of strings $S = \{ \text{bear, bell, bid, bull, buy, sell, stock, stop} \}$

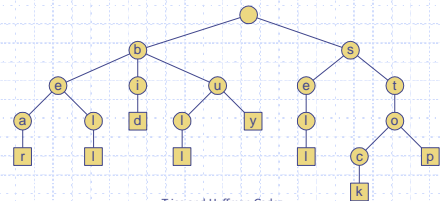


Tries and Huffman Codes

3

Standard Trie

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



Tries and Huffman Codes

4

Application:

- ◆ use a trie to perform a special type of pattern matching: word matching.
- ◆ differ from standard pattern matching since the pattern can not match with an arbitrary substring of the text, but only one of its words.
- ◆ suitable for applications where a series of queries is performed on a fixed text.

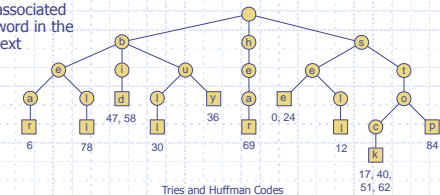
Tries and Huffman Codes

5

Word Matching with a Trie

- ◆ We insert the words of the text into a trie
- ◆ Each leaf stores the occurrences of the associated word in the text

s	e	e	a	b	e	a	r	?	s	e	l	l	s	t	o	c	k	!	!					
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	
s	e	e	a	b	u	l	l	?	b	u	y	s	t	o	c	k	!	!						
24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46		
b	i	d	s	t	o	c	k	!	b	i	d	s	t	o	c	k	!	!						
47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68			
h	e	a	r	!	t	h	e	!	b	e	l	l	?	s	t	o	p	!						
69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88					



Tries and Huffman Codes

6

Compressed Tries:

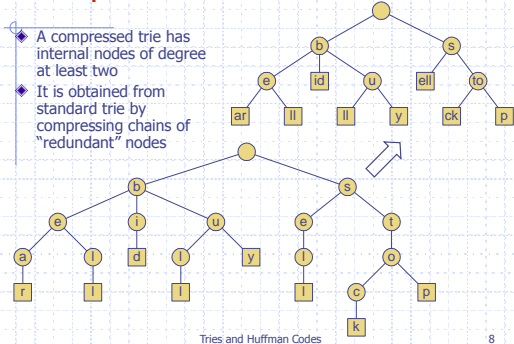
- ◆ An internal node v of T is redundant if v has one child and is not the root.
- ◆ A chain of redundant nodes can be compressed by replacing the chain with a single node with the concatenation of the labels of nodes in the chain.

Tries and Huffman Codes

7

Compressed Trie

- ◆ A compressed trie has internal nodes of degree at least two
- ◆ It is obtained from standard trie by compressing chains of "redundant" nodes



Tries and Huffman Codes

8

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
 - S is an array of strings $S[0], \dots, S[s-1]$
 - Instead of storing a node label X explicitly, we represent it implicitly by a triplet of integers (i, j, k) , such that $X = s[i]j..k$.

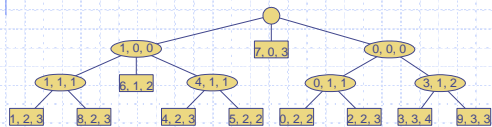
$S[0] =$	0 1 2 3 4	$S[4] =$	0 1 2 3	$S[7] =$	0 1 2 3
$S[1] =$	s e l e	$S[5] =$	b u i l l	$S[8] =$	h e l a r
$S[2] =$	b e l a r	$S[6] =$	b u l y	$S[9] =$	b e l l l
$S[3] =$	s e l l l		b l i d		s t i o p
	s t o c k				

Tries and Huffman Codes

9

Compact Representation

$S[0] =$	0 1 2 3 4	$S[4] =$	0 1 2 3	$S[7] =$	0 1 2 3
$S[1] =$	s e l e	$S[5] =$	b u l l	$S[8] =$	h e l a r
$S[2] =$	b e l a r	$S[6] =$	b u l y	$S[9] =$	b e l l l
$S[3] =$	s e l l l		b l i d		s t i o p



Tries and Huffman Codes

10

Huffman Encoding

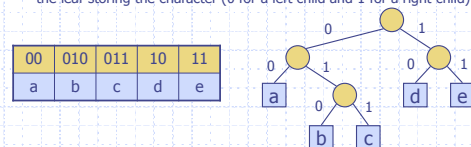
- ◆ Compression
 - Typically, in files and messages,
 - ◆ Each character requires 1 byte or 8 bits (fixed length coding)
 - ◆ Already wasting 1 bit for most purposes!
- ◆ Question
 - What's the smallest number of bits that can be used to store an arbitrary piece of text?
- ◆ Idea
 - Find the frequency of occurrence of each character
 - Encode Frequent characters **short bit strings**
 - Rarer characters **longer bit strings**
 - Variable length coding

Tries and Huffman Codes

11

Encoding Trie

- ◆ 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)



Tries and Huffman Codes

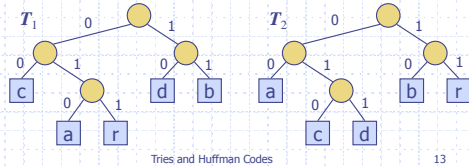
12

Encoding Trie

- Given a text string X , we want to find a prefix code for the characters of X that yields a small encoding for X

Example

- $X = \text{abracadabra}$
- T_1 encodes X into 29 bits
- T_2 encodes X into 24 bits

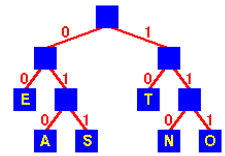


13

Huffman Encoding

Encoding

- Use a tree
- Encode by following tree to leaf
- eg
 - E is 00
 - S is 011
- Frequent characters
 - E, T 2 bit encodings
- Others
 - A, S, N, O 3 bit encodings



Tries and Huffman Codes

14

Huffman Encoding

Greedy Approach

- Sort characters by frequency
- Form two lowest weight nodes into a sub-tree
 - Sub-tree weight = sum of weights of nodes
- Move new tree to correct place

Tries and Huffman Codes

15

Huffman's Algorithm

- Given a string X , Huffman's algorithm construct a prefix code the minimizes the size of the encoding of X
- A heap-based priority queue is used as an auxiliary structure

```

function HuffmanEncoding( $X$ )
  Input string  $X$  of size  $n$ 
  Output optimal encoding trie for  $X$ 
   $C \leftarrow \text{distinctCharacters}(X)$ ;
   $\text{computeFrequencies}(C, X)$ ;
   $Q \leftarrow$  new empty heap;
  for ( $c \in C$ ) {
     $T \leftarrow$  (new single-node tree storing  $c$ );
     $Q.\text{insert}(\text{getFrequency}(c), T)$ ;
  }
  while ( $Q.\text{size}() > 1$ ) {
     $f_1 \leftarrow Q.\text{minKey}()$ ;
     $T_1 \leftarrow Q.\text{removeMin}()$ ;
     $f_2 \leftarrow Q.\text{minKey}()$ ;
     $T_2 \leftarrow Q.\text{removeMin}()$ ;
     $T \leftarrow \text{join}(T_1, T_2)$ ;
     $Q.\text{insert}(f_1 + f_2, T)$ ;
  }
  return  $Q.\text{removeMin}()$ 
    
```

Tries and Huffman Codes

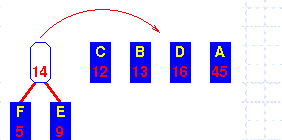
16

Huffman Encoding - Operation

Initial sequence
Sorted by frequency



Combine lowest two
into sub-tree



Move it to correct
place

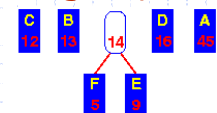


Tries and Huffman Codes

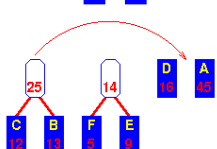
17

Huffman Encoding - Operation

After shifting sub-tree
to its correct place ...



Combine next lowest
pair



Move sub-tree to
correct place

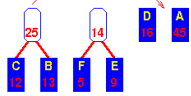


Tries and Huffman Codes

18

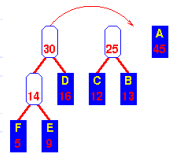
Huffman Encoding - Operation

Move the new tree to the correct place ...



Now the lowest two are the "14" sub-tree and D

Combine and move to correct place

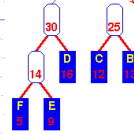


Tries and Huffman Codes

19

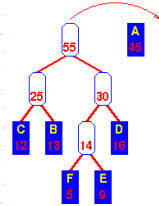
Huffman Encoding - Operation

Move the new tree to the correct place ...



Now the lowest two are the "25" and "30" trees

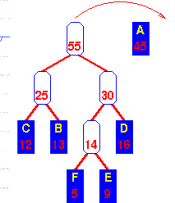
Combine and move to correct place



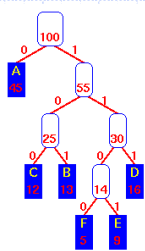
Tries and Huffman Codes

20

Huffman Encoding - Operation



Combine last two trees



Tries and Huffman Codes

21

Example

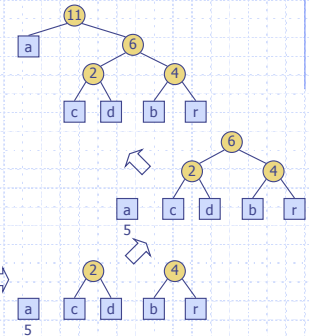
$X = \text{abracadabra}$
Frequencies

a	b	c	d	r
5	2	1	1	2

a	b	c	d	r
5	2	1	1	2

a	b	c	d	r
5	2	1	1	2

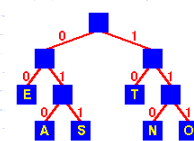
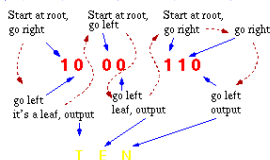
a	b	c	d	r
5	2	1	1	2



Tries and Huffman Codes

22

Huffman Encoding - Decoding



Tries and Huffman Codes

23

Huffman Encoding - Time Complexity

- ◆ Sort keys $O(n \log n)$
- ◆ Repeat n times
 - Form new sub-tree $O(1)$
 - Move sub-tree (binary search) $O(\log n)$
 - Total $O(n \log n)$
- ◆ Overall $O(n \log n)$

Tries and Huffman Codes

24