

Tries & Huffman Codes

CIS2520 - Data Structures

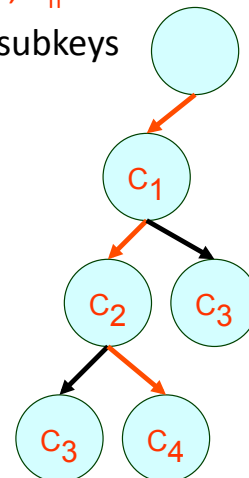
Fall 2011, LAB 9

Tao Xu

1

Trie - Indexed Search Tree

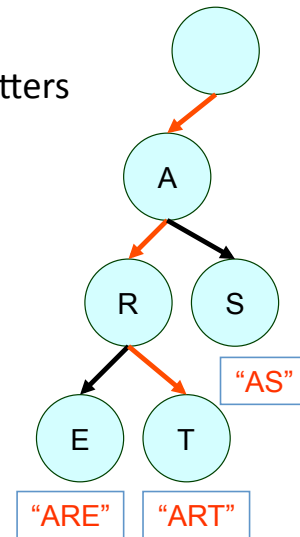
- A trie (retrieval) is a special case of tree
- Applicable when
 - Key **C** can be decomposed into subkeys C_1, \dots, C_n
 - Redundancy (same prefix) exists between subkeys
- Approach
 - Store subkey at each node
 - Path through trie yields full key



2

Tries

- Useful for string matching
 - String decomposed into sequence of letters
 - E.g.,
 - “ART” \Rightarrow “A” “R” “T”
- Can be very fast
 - Less overhead than hashing
- May reduce memory
 - Exploiting redundancy



3

Types of Tries

- Standard
 - Single character per node
- Compressed
 - Eliminating chains of nodes
- Compact
 - Stores indices into original string(s)

4

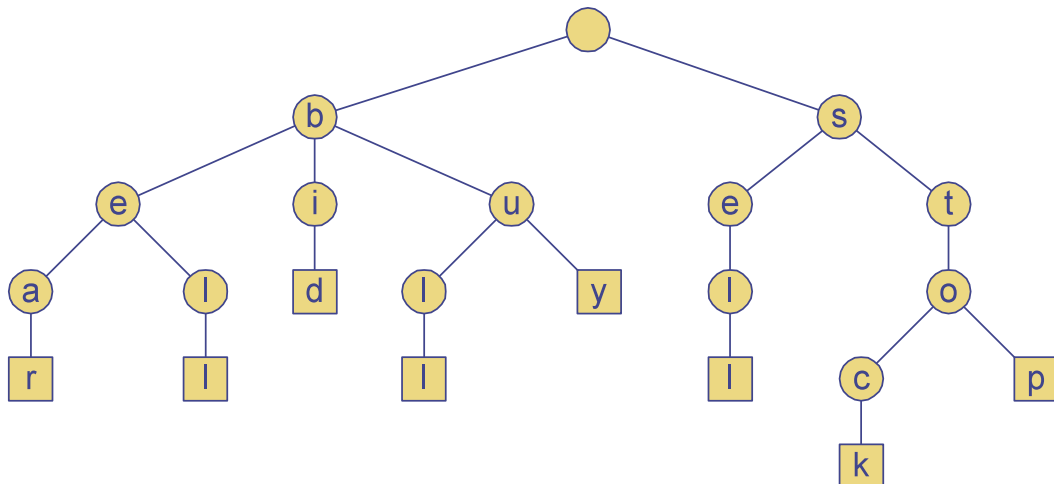
Standard Tries

- A standard trie for a set of strings $S=\{S_1, S_2, \dots\}$
 - Each node (except root) is labeled with a character
 - Children of every node are ordered (alphabetically)
 - Paths from root to leaves yield all strings in S

5

Standard Trie Example

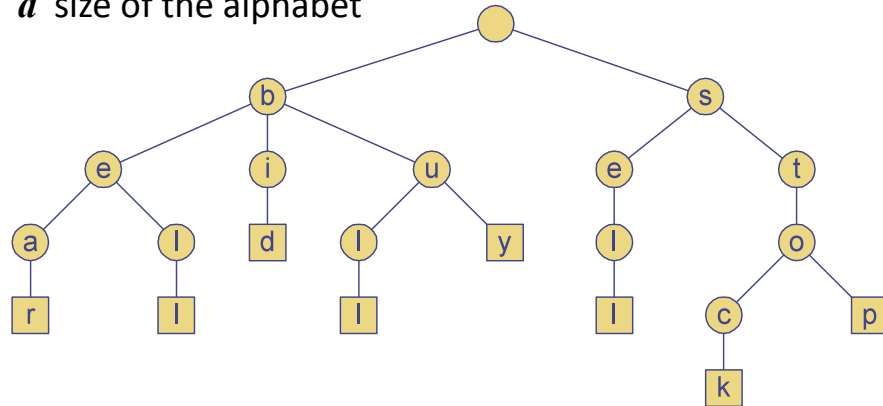
- For strings
 - { bear, bell, bid, bull, buy, sell, stock, stop }



6

Standard Trie (2)

- 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

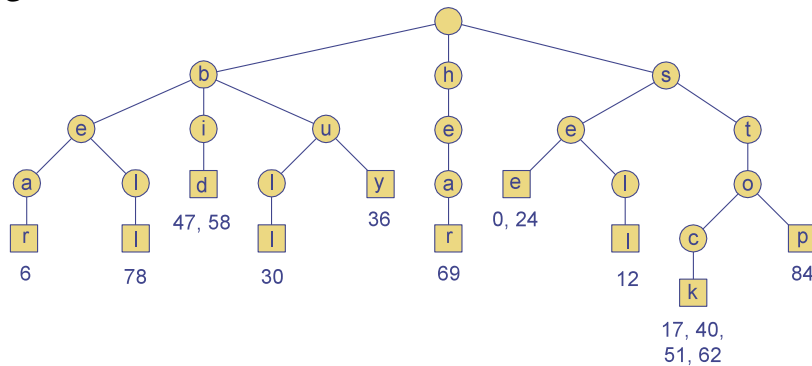


7

Word Matching Trie

- Insert words into trie
- Each leaf stores occurrences of word in the text

| | | | | | | | | | | | | | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| s | e | e | | a | | b | e | a | r | ? | | s | e | i | | 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 | i | 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 | i | 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 | | | | |



8

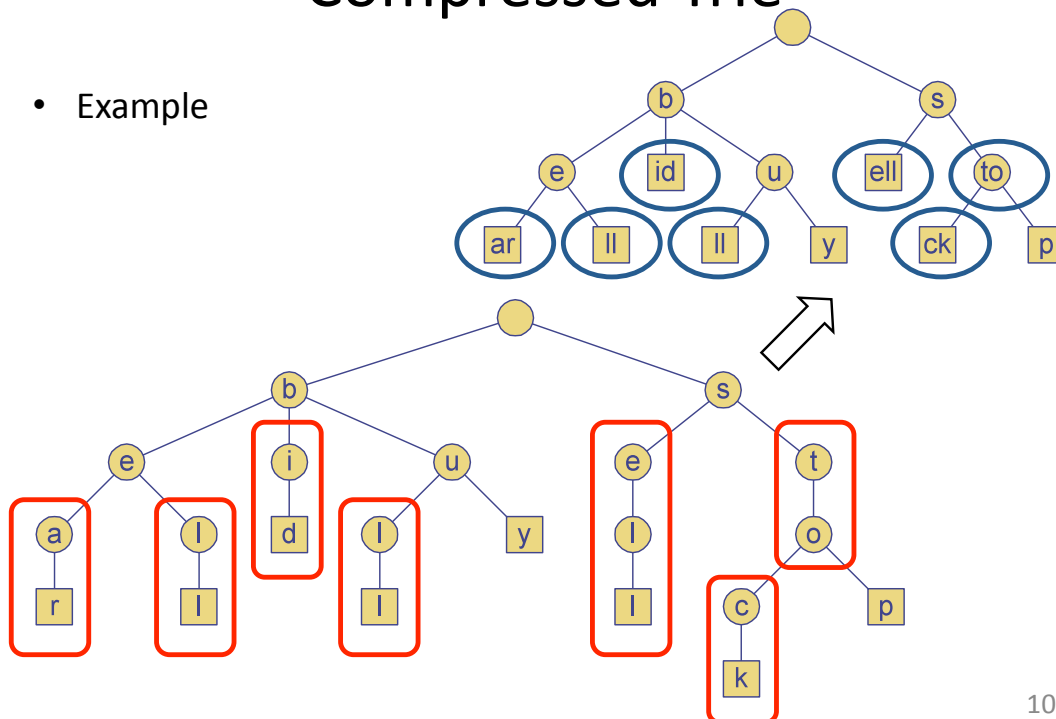
Compressed Trie

- Observation
 - Internal node v of T is redundant if v has one child and is not the root
- Approach
 - A chain of redundant nodes can be compressed
 - Replace chain with single node
 - Include concatenation of labels from chain
- Result
 - Internal nodes have at least 2 children
 - Some nodes have multiple characters

9

Compressed Trie

- Example



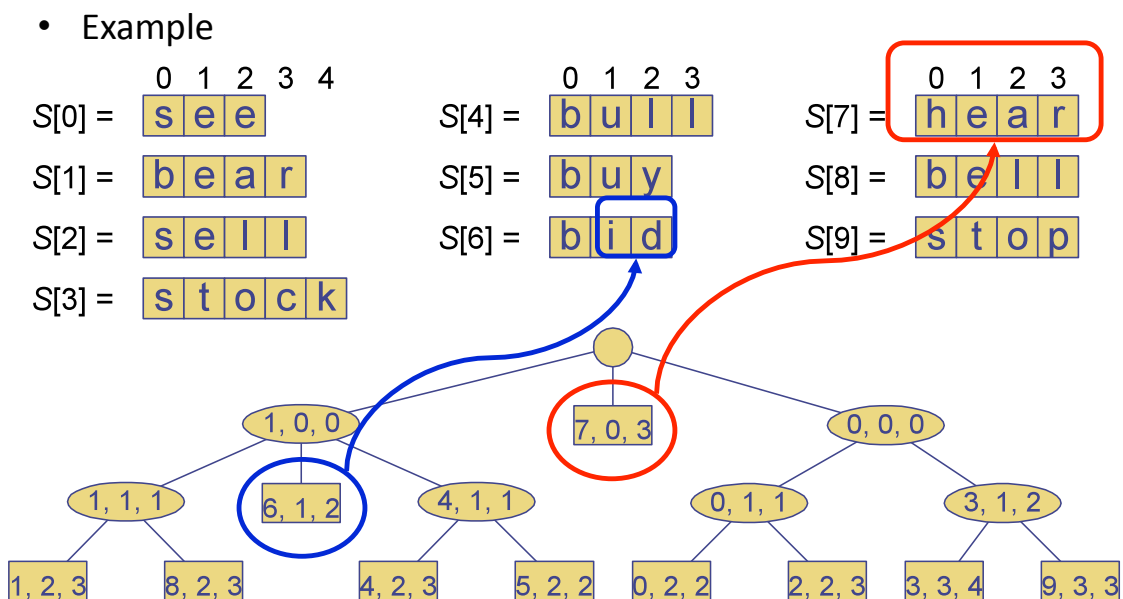
10

Compact Tries

- Compact representation of a **compressed trie**
- Approach
 - For an array of strings $S = S[0], \dots S[s-1]$
 - Store ranges of indices at each node
 - Instead of substring
 - Represent as a triplet of integers (i, j, k)
 - Such that $X = s[i][j..k]$
 - Example: $S[0] = \text{"abcd"}$, $(0,1,2) = \text{"bc"}$
- Properties
 - Uses $O(s)$ space, where $s = \#$ of strings in the array
 - Serves as an auxiliary index structure

11

Compact Representation



12

Tries and its Applications

- Web search engine indexing
- Spelling check
- Computational biology
- ...

13

Huffman Encoding

- Compression
 - Typically, in files and messages,
 - Each character requires 1 byte or 8 bits
 - Already wasting 1 bit for most purposes!
- Question
 - What's the minimum 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**
 - Rare characters -- **long bit strings**

14

Encoding – A Simple Example

- A message of 5 symbols:

[▶♣♠♥☺♣☼☺]

- Fix length code (e.g, ASCII)
5 symbols → at least 3 bits each
Total Length = $5 \times 3 = 15$ bits

| | |
|---|-----|
| ▶ | 000 |
| ♣ | 001 |
| ☺ | 010 |
| ♠ | 011 |
| ☼ | 100 |

- Huffman code
Length varies depending on frequency
Total Length = $3 \times 2 + 3 \times 2 + 2 \times 2 + 1 \times 2 + 1 \times 2 = 24$ bits

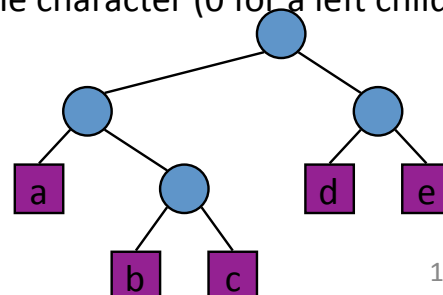
| Symbol | Freq. | Code |
|--------|-------|------|
| ▶ | 3 | 00 |
| ♣ | 3 | 01 |
| ☺ | 2 | 10 |
| ♠ | 1 | 110 |
| ☼ | 1 | 111 |

15

Encoding Trie

- A **code** is a mapping 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)

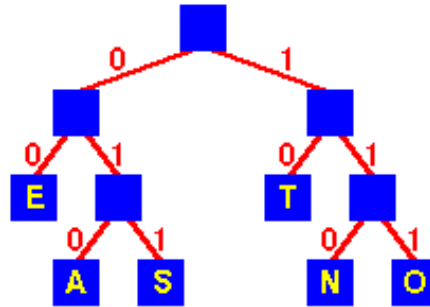
| | | | | |
|----|-----|-----|----|----|
| 00 | 010 | 011 | 10 | 11 |
| a | b | c | d | e |



16

Huffman Encoding

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



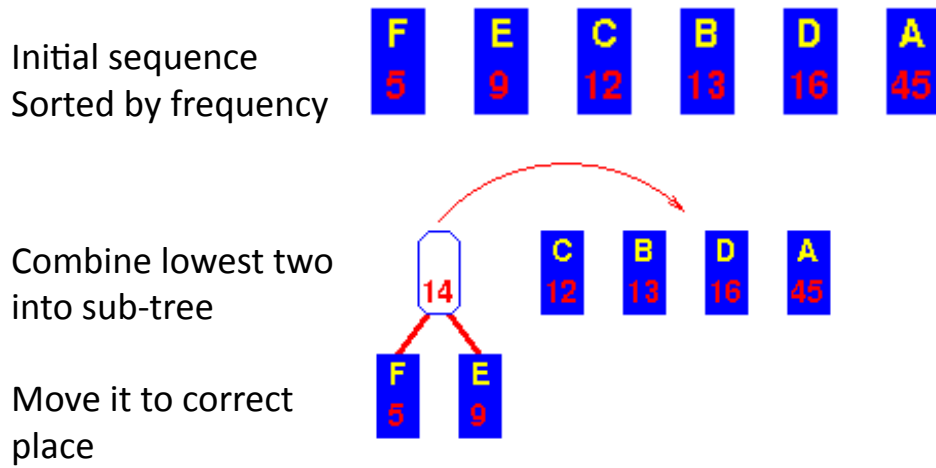
17

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
- Bottom-up (optimal)
 - Top-down (Shannon coding, suboptimal)
- Suboptimal to encode block data (e.g., 'cat')

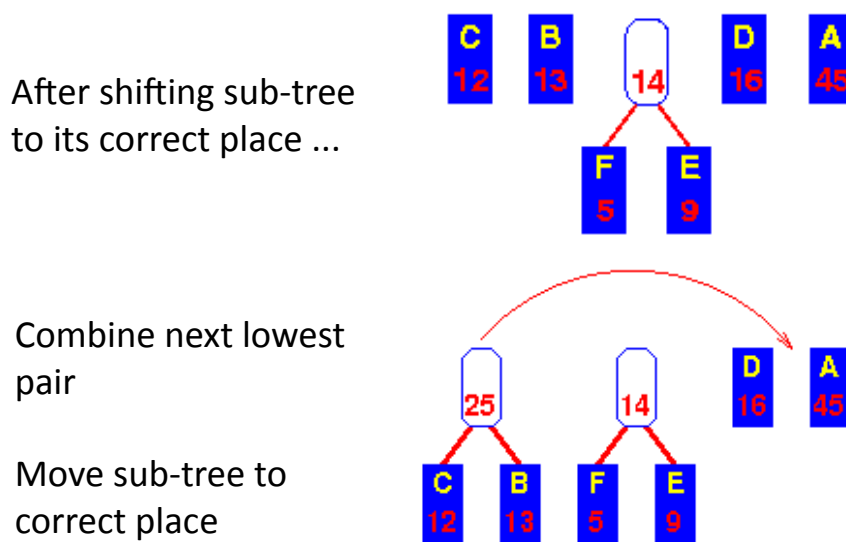
18

Huffman Encoding - Operation



19

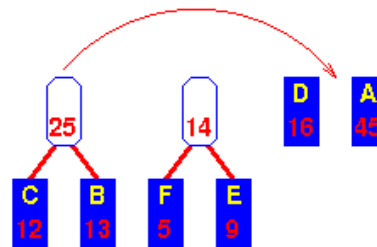
Huffman Encoding - Operation



20

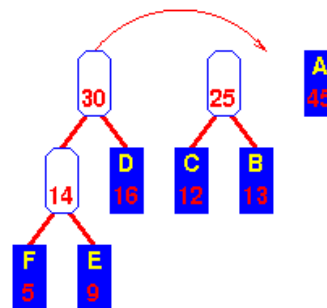
Huffman Encoding - Operation

Move the new tree to the correct place ...



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

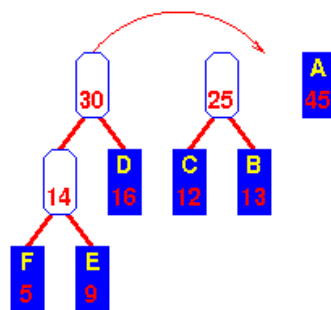
Combine and move to correct place



21

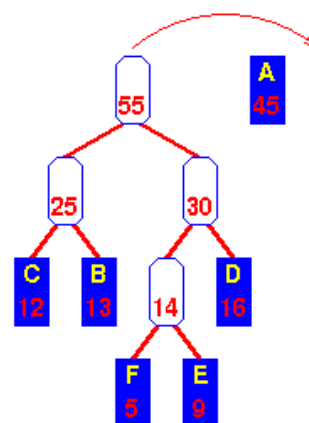
Huffman Encoding - Operation

Move the new tree to the correct place ...



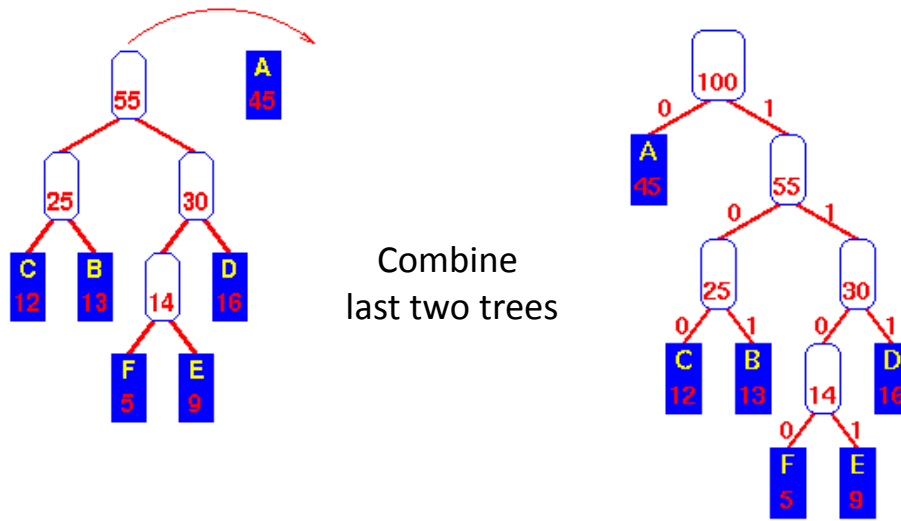
Now the lowest two are the the "25" and "30" trees

Combine and move to correct place



22

Huffman Encoding - Operation



23

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  $X$ 
   $C = \text{distinctCharacters}(X)$ 
   $\text{computeFrequencies}(C, X)$ 
   $Q = \text{new empty heap}$ 
  for all  $c$  in  $C$ 
     $T = (\text{new single-node tree storing } c)$ 
     $Q.\text{insert}(\text{getFrequency}(c), T)$ 
  while  $Q.\text{size}() > 1$ 
     $f_1 = Q.\text{minKey}()$ 
     $T_1 = Q.\text{removeMin}()$ 
     $f_2 = Q.\text{minKey}()$ 
     $T_2 = Q.\text{removeMin}()$ 
     $T = \text{join}(T_1, T_2)$ 
     $Q.\text{insert}(f_1 + f_2, T)$ 
  return  $Q.\text{removeMin}()$ 
  
```

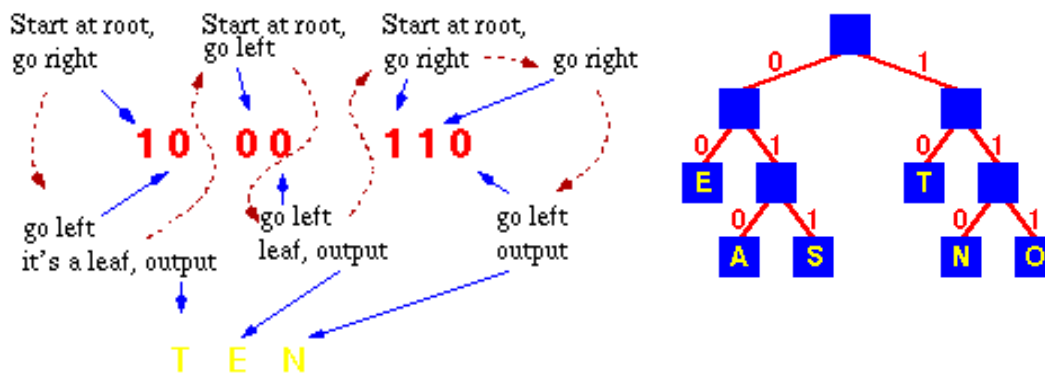
24

Huffman Encoding - Time Complexity

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

25

Huffman Encoding - Decoding



26