# UNIT-4
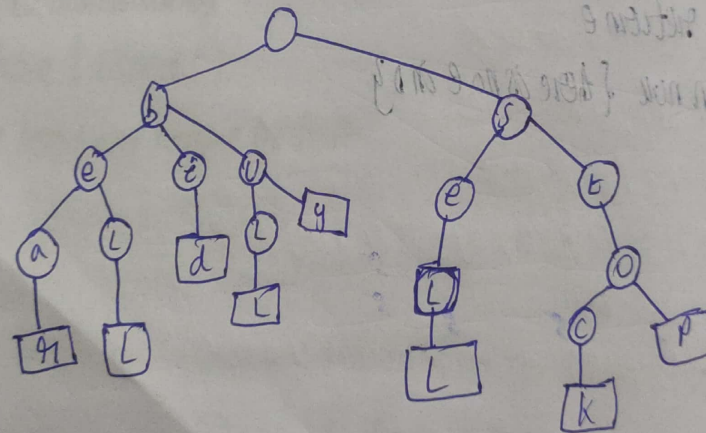
Text processing : String Operations ; Brute-force Pattern matching ;
The Boyer moore Algorithm , The knuth - morris Pratt Algorithm
standard tries , compressed tries , suffix tries , Huffman
coding algorithm , the longest subsequence Problem (LCS),
Applying dynamic programming to the LCS Problem.

TRIE:
A data structure for storing a set of strings (Name from the
word "retrieval"):

# PROPERTIES OF TRIE:

(1) Multiway tree

(2) Each node has from 1 to d children.

(3) Each edge of the tree is labeled with a Character.

(4) Each leaf node corresponds to the Stored string, which is a concatenation of characters on a path from the root to this node.
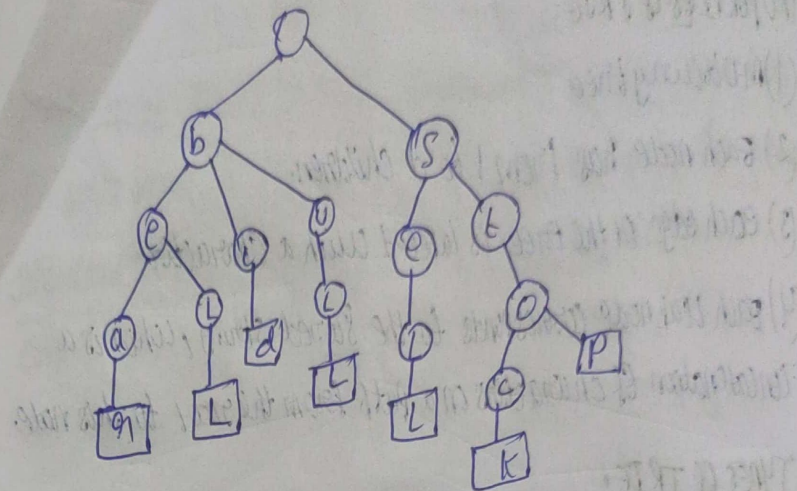
## TYPES OF TRIE:

(1) Standard trie

(2) Compressed trie

(3) Suffix trie

## STANDARD TRIE:

1. Each node but the root is labelled with a character

2. The children of a node are alphabetically ordered.

3. The paths from the external nodes to the root yield to the strings of S.

4. Example:- standard trie for the set of strings S

$S = \{ bear, bell, bid, bull, buy, sell, stock, stop \}$

→ Standard trie uses $O(n)$ space. Operations (find, insert, remove) take time $O(dm)$ each, where:
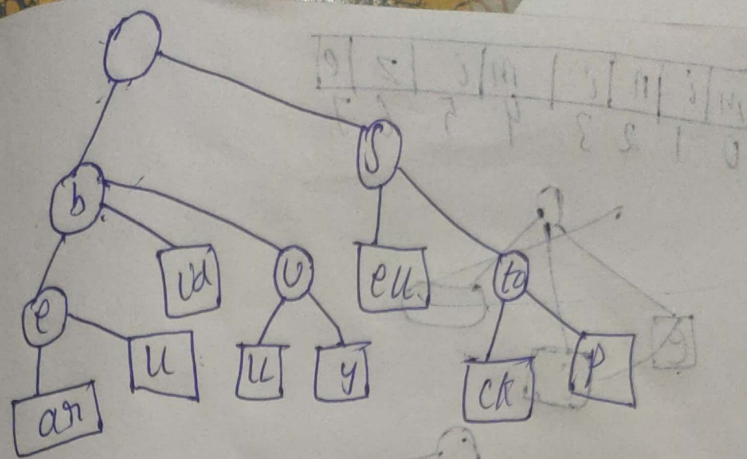
        $n$ ~ total size of strings ins

        $d$ – alphabet size

        $m$ ~ size of the string parameter of the operation

Compressed tries:

(1) Tree with nodes of degree at least 2

(2) Obtained from standard trie by compressing the chains of redundant blocks

A compressed trie can be stored in space $O(s)$ where $s = |r|$ by using $O(1)$ size nodes.

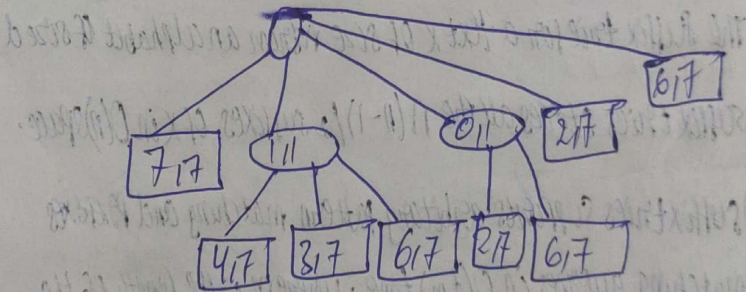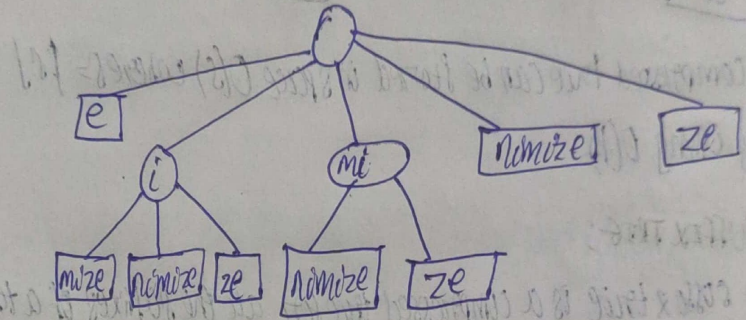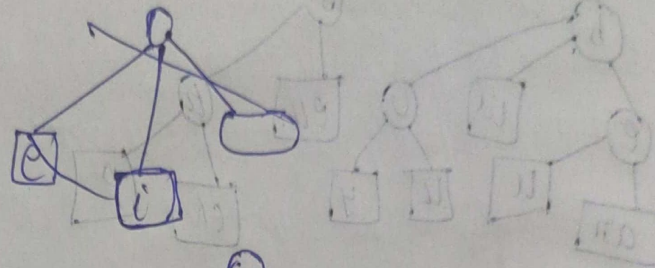## SUFFIX TRIE:

1) A suffix trie is a compressed trie for all the suffixes of a text.

2) The suffix trie for a text x of size n from an alphabet of size d.

3) Suffix tries stores all the $n(n-1)/2$ suffixes of x in $O(n)$ space.

4) Suffix tries supports arbitrary pattern matching and prefix matching queries in $O(dm)$ time, where m is the length of the pattern.

5) Suffix tries can be constructed in $O(dn)$ time.
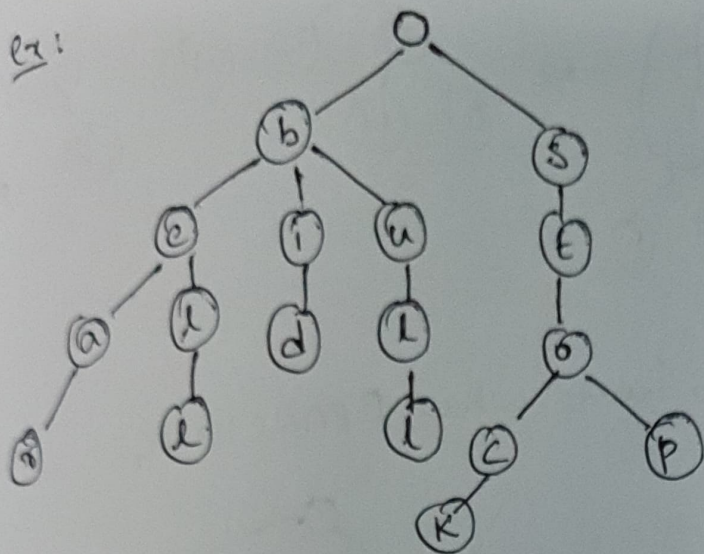
## Applications:

(1) word matching

(2) prefix matching

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

## Standard Trie:

In standard trie no word in s should be the prefix of other.

ex:



→ class Node {

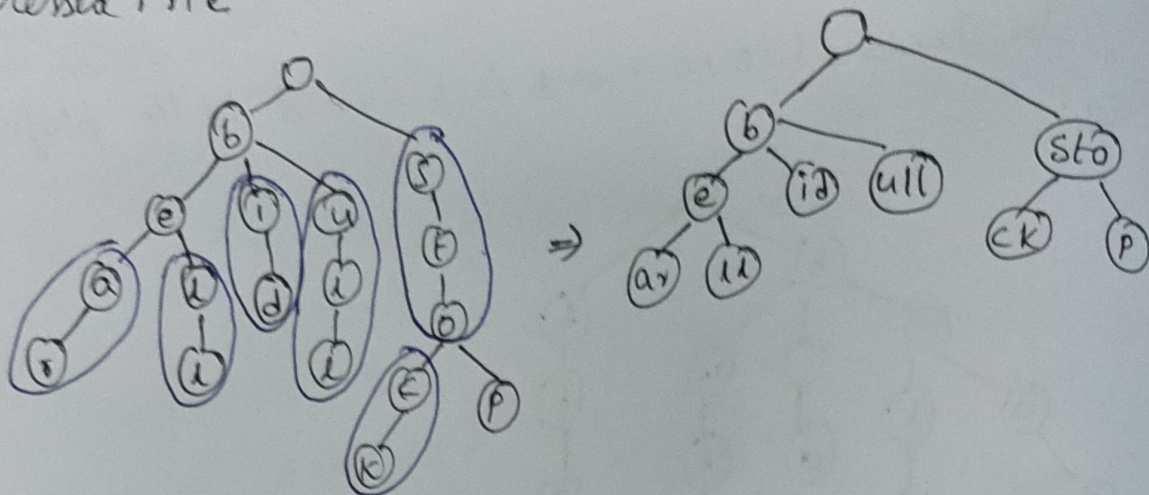        Node [] children = new node [26];

        // to check for end of string

        boolean iswordend;

}

→ It is an ordered tree like data structure

→ Each node (except the root node) is a standard trie is labeled with a character.

→ The children of a node are in alphabetical order.

→ Each node or branch represents a possible character of keys or words.

→ Each node or branch may have multiple branches.

→ The last node of every key or word is used to mark the end of word or node.

# Compressed Trie



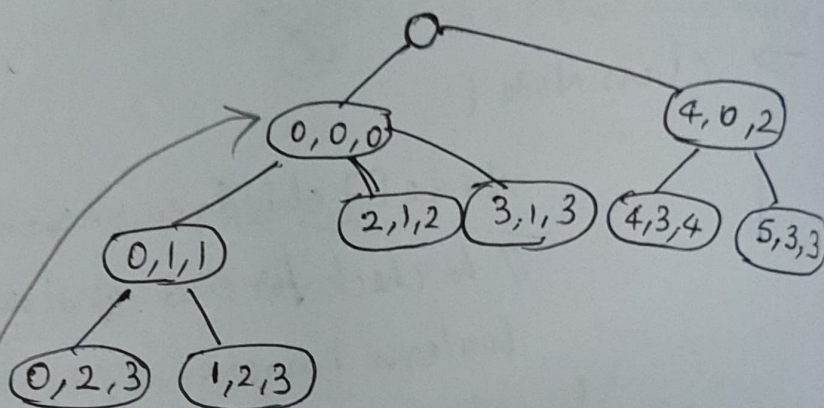NOTE : only add a new node when BRANCHING

$$i \quad 0 1 2 3 4$$

$S[0] = bear$

$s[1] = bell$

$S[2] = bid$

$S[3] = bull$

$s[4] \quad stock$

$s[5] = stop$

Node $(i, j, k)$

$i$ – index of $S$

$j$ – start

$k$ – end

# Suffix Node Trie:

Ex: Minimize

Suffixes:

$$\left.\begin{array}{l} e \\ ze \\ ize \\ mize \\ imize \\ nimize \\ inimize \\ Minimize \end{array}\right\} - S$$



→ Now Compress



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

(start, end)