

Boyer-Moore

- search from last index of pattern
- keep a table of indexes in the pattern
- if mismatch, turn it into a match:

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
G C A A T G C C T A T G T G A C C
T A T G T G

```

- if the last character mismatches, move entire pattern:

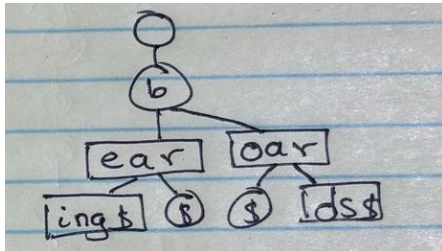
```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
G C A A T G C C T A T G T G A C C
O T A T G T G

```

PATRICIA

- assume a node is redundant if it has one child and it's not the root.
- chain redundant nodes:



Bloom Filter

- pass search key through k hash functions.
- check if result of hash functions in your table is 1, then the searched key exists in your set.
- if any bit from the hashed result is 0---does not exist.
- collisions can occur, resulting in false positives.

n : # of items
 f : false positive rate
 m : # of bits
 k : # of hash functions

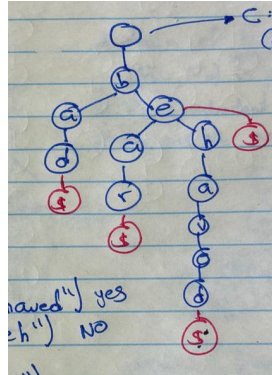
$$k_{opt} = \frac{m}{n} \ln 2$$

$$f_{opt} = \left(\frac{1}{2}\right)^{k_{opt}}$$

$$f \approx \left(1 - e^{-\frac{k}{m/n}}\right)^k$$

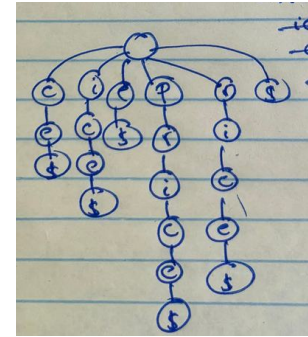
Tries

- structure for patterns using ordered alphabets
- optional '\$' signifying pattern end (kinda like "\0")
- $O(n)$ space complexity, n is the total length of all strings
- $\sim O(L)$ time complexity, L is the average length of a string



Suffix Trie

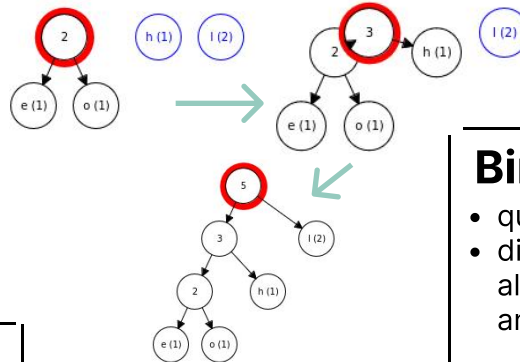
- form all possible suffixes from the string, terminate each with '\$'
- if a path ends with '\$', it's a *suffix*. any path is a *substring*.



- to find the *longest repeated substring*, find the deepest node with more than one child.

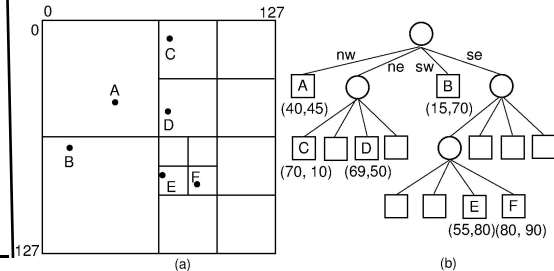
Huffman Coding

- repeatedly combine the 2 *least frequent* nodes into a subtree



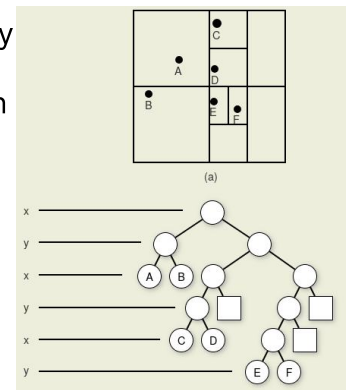
Point-Region Quadtree

- each node has exactly 4 children: nw, ne, sw, se



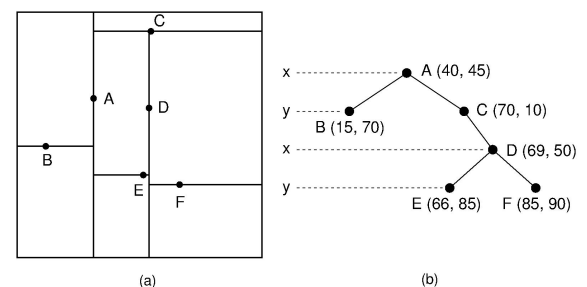
Bintree

- quadtree but binary
- discriminator alternates between x and y coordinates



K-D Tree

- bintree but adaptive



Count min-sketch

- Like a bloom filter, but instead of just storing either 1 or 0 in the table, we store multiple bits representing the integer count of the cell.

ϵ : band of overestimate
 δ : failure probability

$$k = \ln(1/\delta)$$

$$m = \frac{e}{\epsilon}$$