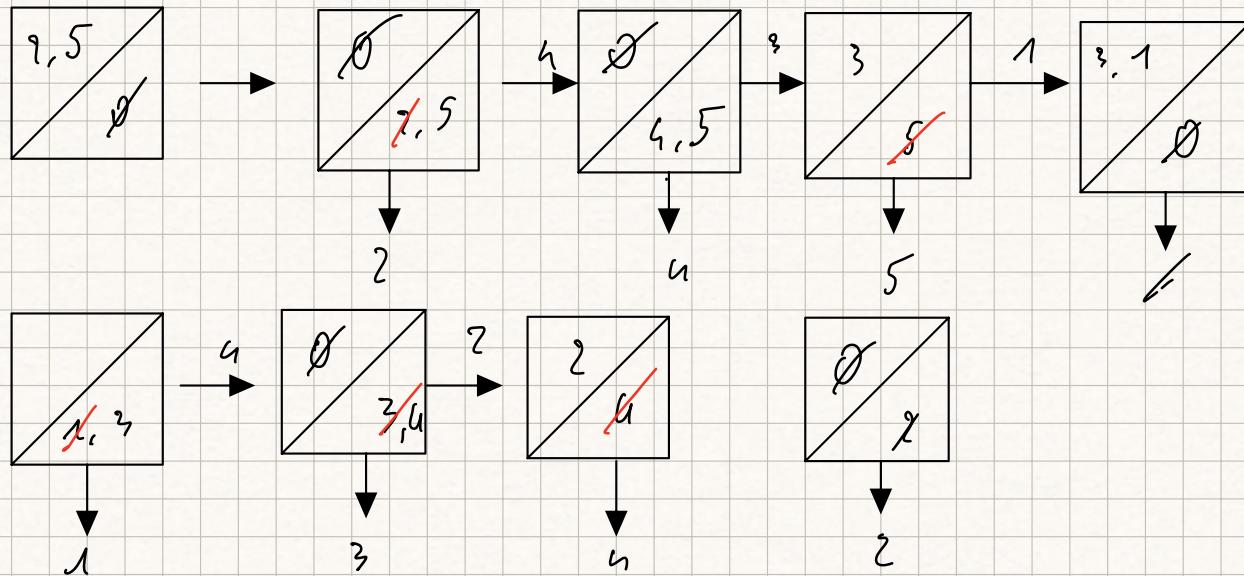


Simulate the algorithm SnowPlow over the sequence 2,5,4,3,1,4,2, and show which sorted blocks it forms with a memory of size M=2.



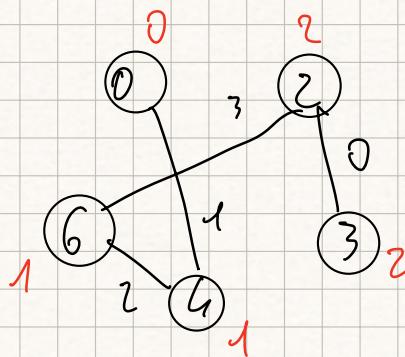
Given the ordered set of strings $S = \{AA, AC, BB, CC\}$. Compute the Minimal Ordered Perfect Hash for S by assuming the following two hash functions:

$$h_1(xy) = x+y \bmod 7 \text{ and } h_2(xy) = x+2^*y \bmod 7$$

in which x (resp. y) is the code of the first (resp. second) letter of a string of S , and the codes are: A=1, B=2, and C=3.

As an example, if the string is AC, then $x=1$ and $y=3$.

	h_0	h_1	h_2
AA	0	2	3
AC	1	4	0
BB	2	4	6
CC	3	6	2



$$\begin{aligned} m' &= 7 \\ \# \text{Keys} &= 4 \end{aligned}$$

$m' \geq 4$ ok

$\text{mod } 4$

$$y = \begin{bmatrix} 0 & 0 & 2 & 2 & 1 & 0 & 1 \\ 0 & 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix}$$

Given the sequence a,b,c,d,e,f,g,h,i,l simulate:

- The sampling algorithm for m=2 which knows the sequence length n=10,

assuming probabilities for the parameter p = [0.5, 0.5, 0.5, 1, 1, 0.1, 0.5, 1, 0.1, 1]

- The sampling algorithm for m=2 which does not know the sequence length,

assuming values for the parameter h = [1, 3, 4, 2, 1, 5, 4, 6]

	1	2	3	4	5	6	7	8	9	10
1) S=0	b	c	d	e	h	g	h	i	l	
P = 0.5	0.5	0.5	0.5	1	1	0.1	0.5	1	0.1	1

$$S=0 \quad J=1 \Rightarrow \frac{2-0}{10-1+1} = \frac{1}{5} \times$$

$$\text{checked if } P \leq \frac{m-j}{m-j+1}$$

$$S=0 \quad J=2 \quad \frac{2-0}{10-2+1} = \frac{2}{9} \times$$

$$S=0 \quad J=6 \quad \frac{2-0}{10-6+1} = \frac{2}{5} \quad \checkmark$$

$$S=0 \quad J=3 \quad \frac{2-0}{10-3+1} = \frac{2}{8} \times$$

$$S=1 \quad J=7 \quad \frac{2-1}{10-7+1} = \frac{1}{4} \times$$

$$S=0 \quad J=4 \quad \frac{2-0}{10-4+1} = \frac{2}{7} \times$$

$$S=1 \quad J=8 \quad \frac{2-1}{10-8+1} = \frac{1}{3} \times$$

$$S=0 \quad J=5 \quad \frac{2-0}{10-5+1} = \frac{2}{6} = \frac{1}{3}$$

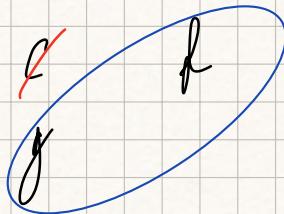
$$S=1 \quad J=9 \quad \frac{2-1}{10-9+1} = \frac{1}{2} \quad \checkmark$$

unknown length

	1	2	3	4	5	6	7	8	9	10
S=0	b	c	d	e	h	g	h	i	l	
P = -	-	-	1	3	6	2	1	5	1	6

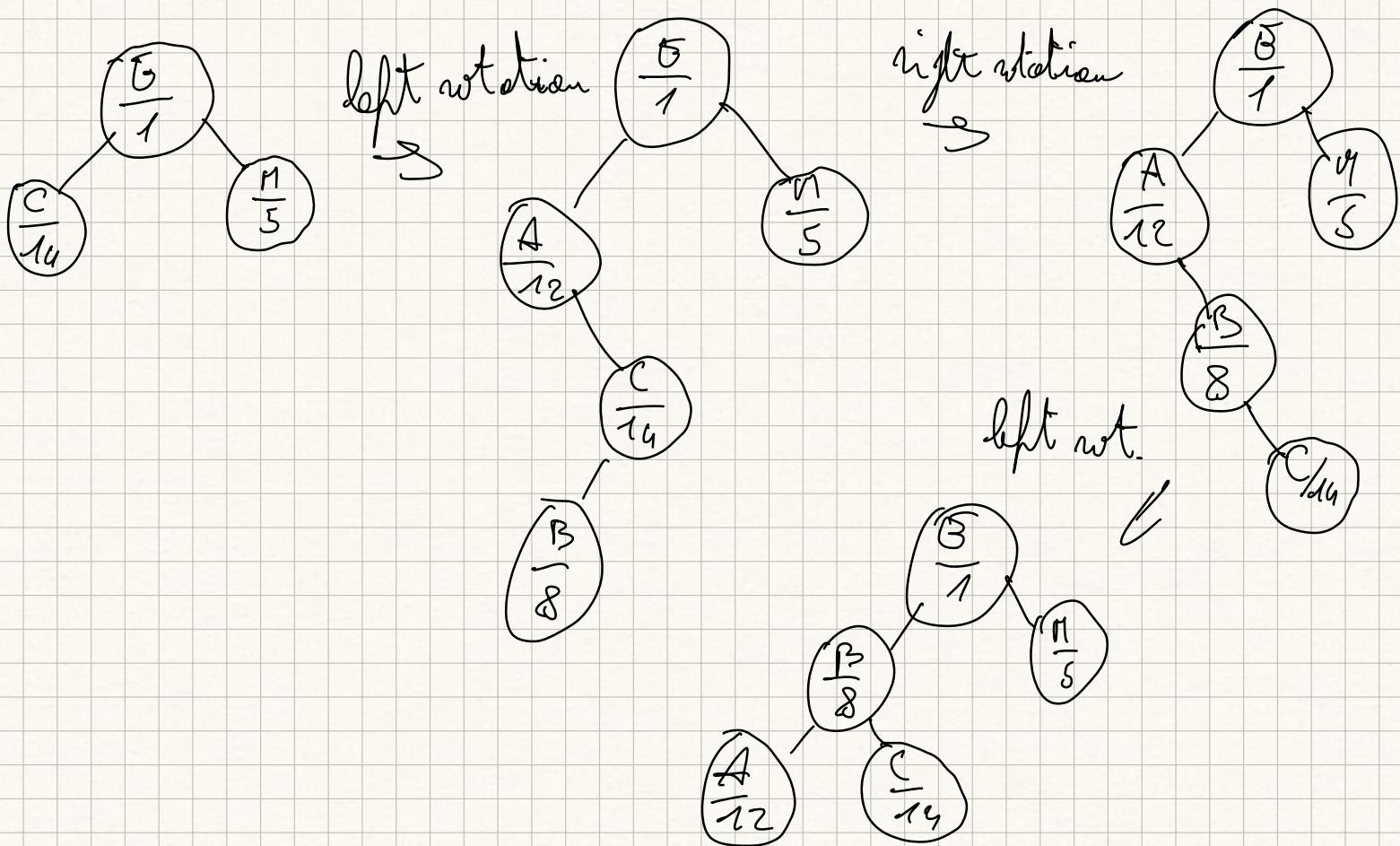
$R =$	<table border="1"> <tr> <td>e</td> <td>b</td> </tr> </table>	e	b
e	b		

initial state

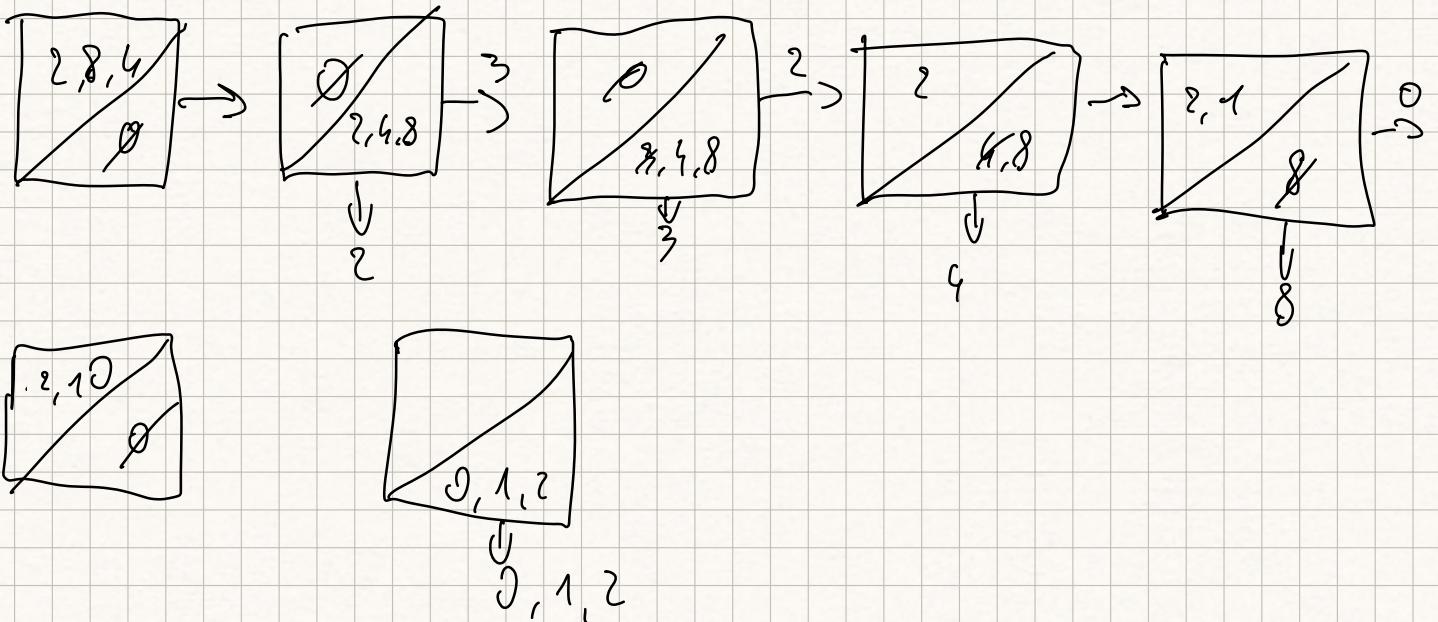


final state

Build a Treap by inserting the following sequence of pairs $\langle \text{key}, \text{priority} \rangle$ and assuming that the MIN priority is in the root (the order among the keys is the alphabetic one):
 $\langle E, 1 \rangle \langle C, 14 \rangle \langle M, 5 \rangle \langle A, 12 \rangle \langle B, 8 \rangle$

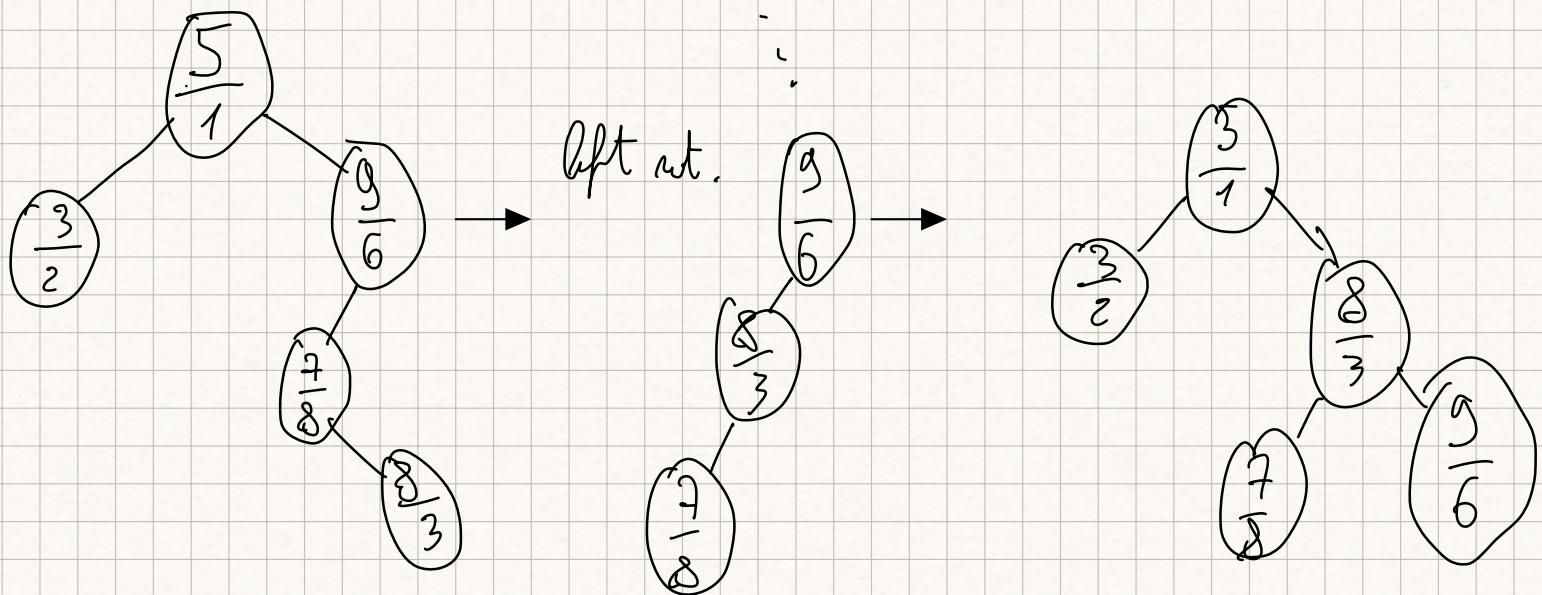


Simulate the algorithm Snow Plow on the sequence: (2, 8, 4, 3, 2, 1, 0) by assuming a memory size $M=3$.



You are given a set of pairs (key, priority) to insert in a Treap according to the following order: (5,1), (9,6), (3,2), (7,8).

- Show the Treap resulting from the insertion of every key above
- Insert then the pair (8,3)



Sort the following strings via Multi-key Quicksort

$S = \{ \text{castro}, \text{abba}, \text{mom}, \text{camel}, \text{astral}, \text{asso} \}$ by using as pivot always the first string of the set to sort.

Pivot = castro $i = 1$

$$P = ABBA$$

$$R < \{ABRA, ASTRA, ASSO\} \quad R_C = \emptyset \quad Q = \{ABRA, ASTRA, ASSO\} \quad Q_{>} = \emptyset$$

$$P = ABRA$$

$$R < \emptyset \quad R_C = \emptyset \quad Q = \{ABRA\} \quad Q_{>} = \{ASTRA, ASSO\}$$

$$P = ASTRA$$

$$R < \emptyset \quad R_C = \emptyset \quad Q = \{ASTRA\} \quad Q_{>} = \emptyset$$

$$P = CASTRO, CANC$$

$$R < \emptyset \quad R_C = \emptyset \quad Q = \{CASTRO, CANC\}$$

$$P = \emptyset$$

$$R < \emptyset \quad R_C = \emptyset \quad Q = \{CASTRO, CANC\} \quad Q_{>} = \emptyset$$

$$P = \emptyset$$

$$R < \emptyset \quad R_C = \emptyset \quad Q = \{CASTRO\} \quad Q_{>} = \emptyset$$

$R > \text{none}$

Given the following set of strings, each consisting of two digits:

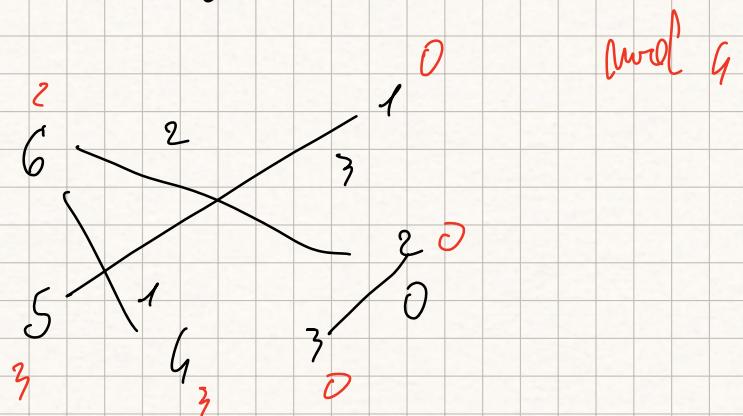
$$S = \{11, 22, 33, 44\}$$

Compute the Minimal Ordered Perfect Hash for S by assuming the following two hash functions: $h_1(xy) = x+y \bmod 7$ and $h_2(xy) = x+2^y \bmod 7$

in which x (resp., y) is the first (resp., second) digit of a string of S.

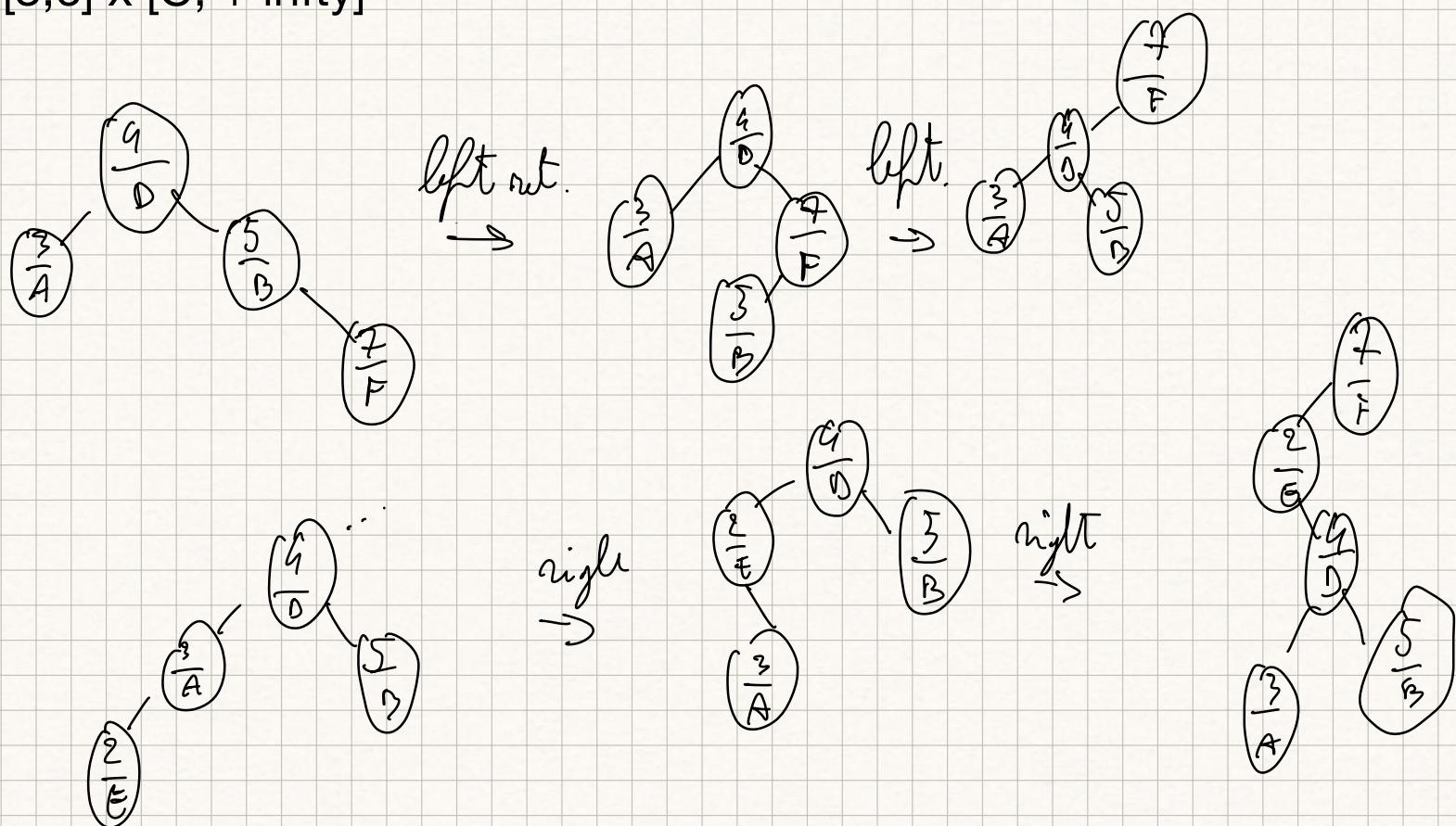
As an example, if the string is 11, then x=1 and y=1.

	h	h_1	h_2	$m' = 7 > m = 4$ ok
11	0	2	3	
22	1	4	6	
33	2	6	2	
44	3	1	5	



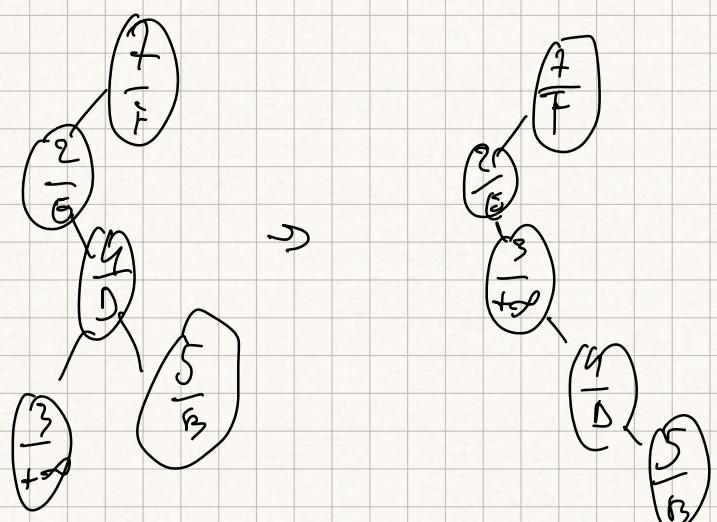
Given the set S of pairs { $\langle 4, D \rangle$, $\langle 5, B \rangle$, $\langle 3, A \rangle$, $\langle 7, F \rangle$, $\langle 2, E \rangle$ }, where the first component is the key and the second component is the priority:

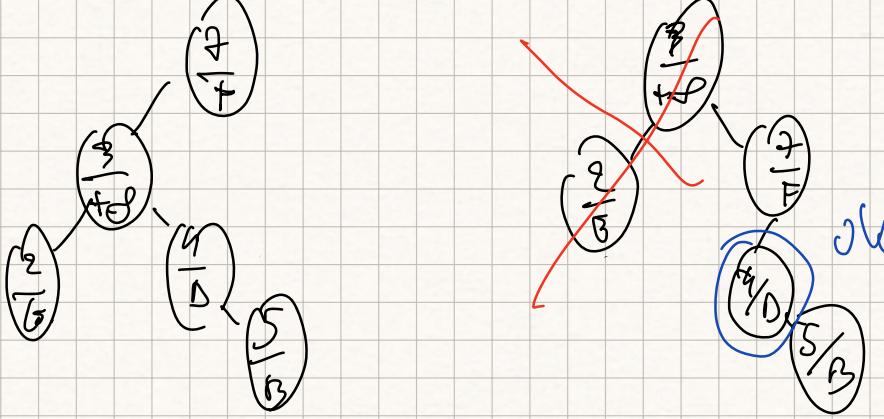
- Build a TREAP data structure by inserting the pairs in that order (you can assume that it is a MAX heap and letters on the Y-axes are sorted from bottom to top as A, B, C,).
- Show and comment how to solve the three-sided range query $[3, 6] \times [C, +\infty]$



2) • Split (3)

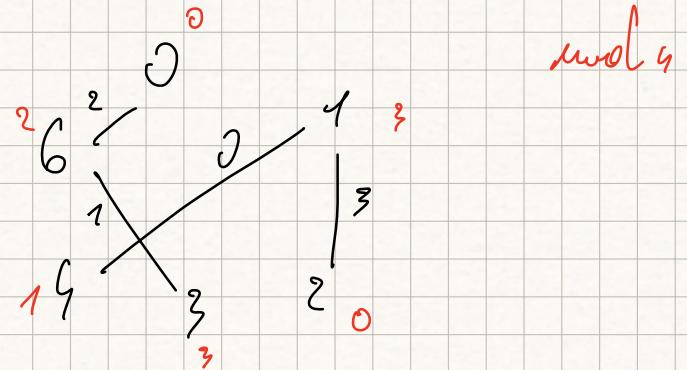
• drop left subtree ('item < 3')





Assume you are given a set of 4 strings $\{aa, ac, bc, cc\}$ and you wish to construct a minimal ordered perfect hash function (MOPHF), where the order is the alphabetic one. Assume that $\text{rank}(x) = 2; 3; 4$ for the characters $x = a; b; c$; respectively. Given a string xy of two characters, we let the two random functions required by the design of MOPHF be $h_1(xy) = (3 * \text{rank}(x) + \text{rank}(y)) \bmod 7$ and $h_2(xy) = (\text{rank}(x) + \text{rank}(y)) \bmod 7$.

l	h_1	h_2
AA	0	1
AC	1	3
BC	2	6
CC	3	2



Given the sequence of 6 items $S = (a, b, c, d, e, f)$ use the random sampling algorithm with known sequence length $n=6$, to compute the $m=2$ extracted items given the random values $p = (1/2, 1/10, 3/4, 3/4, 0, 1)$

$$\begin{aligned} j &= 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \\ S &= a, b, c, d, e, f \\ p &= 1/2 \quad 1/10 \quad 3/4 \quad 3/4 \quad 0 \quad 1 \end{aligned}$$

checked if $P \leq \frac{m-j}{n-j+1}$

$$S=0 \quad j=1 \quad \frac{2-0}{6-1+1} = \frac{1}{3} \times$$

$$S=0 \quad j=2 \quad \frac{2-0}{6-2+1} = \frac{2}{5} > 0.1 \text{ OK}$$

$$S=1 \quad j=3 \quad \frac{2-1}{6-3+1} = \frac{1}{4} \times$$

$$S=1 \quad j=4 \quad \frac{2-1}{6-6+1} = \frac{1}{1} \times$$

$$S=1 \quad j=5 \quad \frac{2-1}{6-5+1} = \frac{1}{2} > 0 \quad \checkmark$$

$$S=1 \quad j=6 \quad \frac{1}{1} = 1 \quad \text{OK}$$

Consider the Snow Plow technique with memory M=2.

- Simulate Snow Plow over the sequence S = (1, 3, 9, 10, 7, 6, 5, 4, 3, 8).

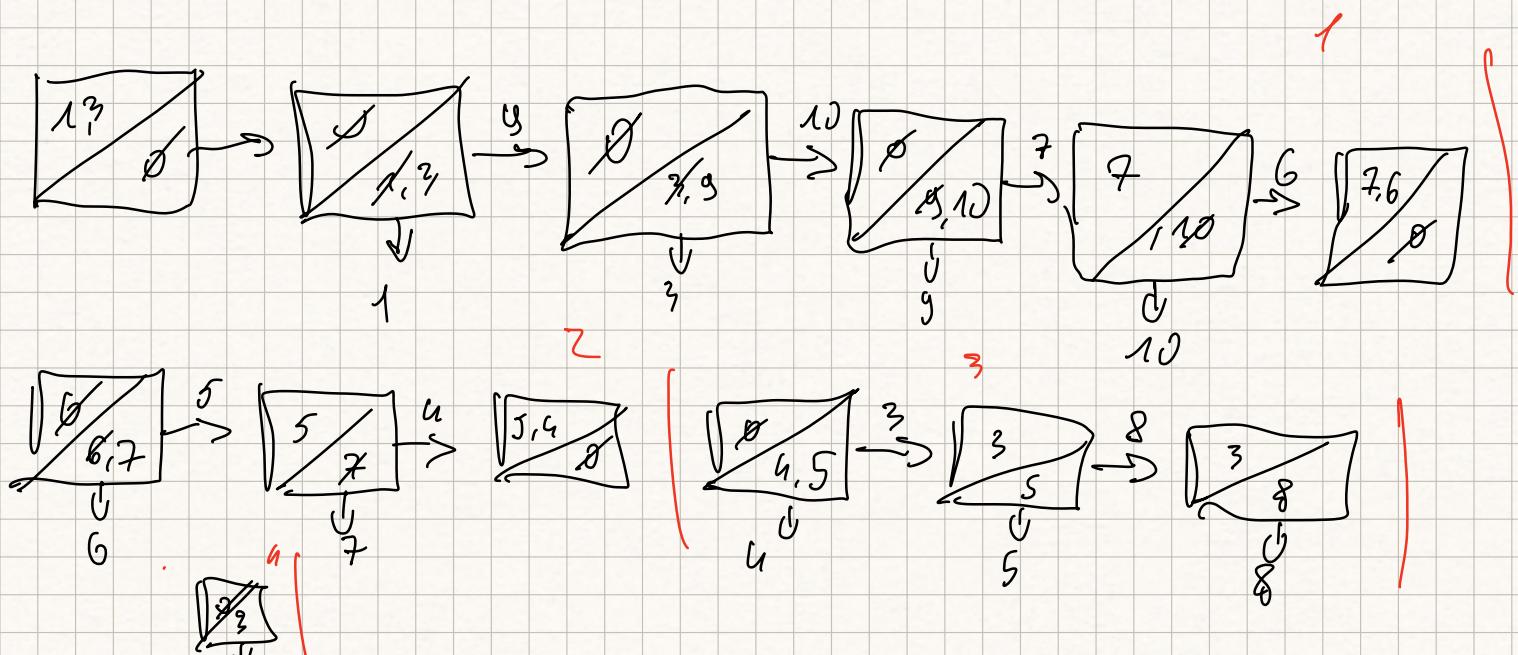
- Provide an example of sequence of length 10 that generates exactly 5 runs,

when the memory has size M=2.

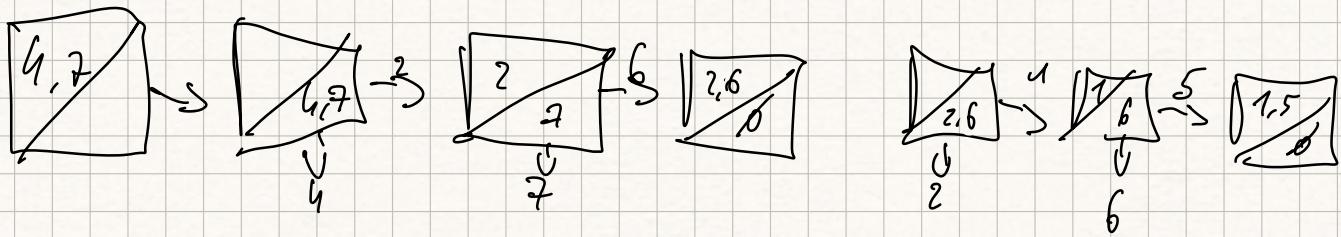
- Show the average length of the run produced by Snow Plow, if we assume

that the probability that an item goes to the Heap is $3/4$ rather than $1/2$.

$$M=2$$



$$S = \{4, 7, 2, 6, 1, 5, 3, 11, 9, 20\}$$



$T_{L_1} = 0$ $T = 4M$ T item read in one phase

$P(\text{item goes to memory}) = 1/4$

$$L_1 = \{1, 2, 4, 6, 8, 10, 15, 18, 20\}$$

$$L_2 = \{2, 3, 7, 8, 10, 18\}$$

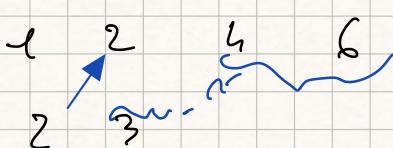
- ① Compute $L_1 \cap L_2$ by mutual partitioning
- ② \rightarrow By 2-level storage approach $B=3$

①

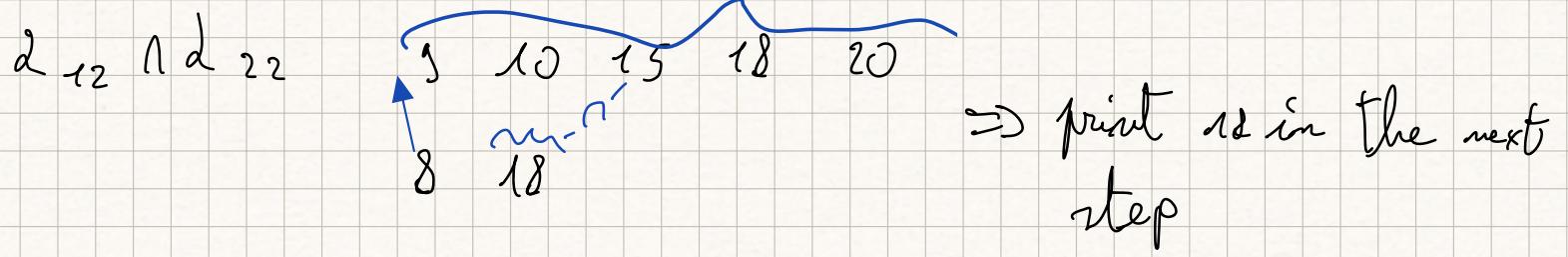
$$L_1 = \{1, 2, 4, 6, 8, 10, 15, 18, 20\}$$

$$L_2 = \{2, 3, 7, 8, 10, 18\}$$

$$L_{11} \cap L_{21}$$



Point 2



②

$$d_1 = \boxed{1 \ 2 \ 6} \quad B_1 \quad \boxed{6 \ 9 \ 10} \quad B_2 \quad \boxed{15 \ 18 \ 20} \quad B_3$$

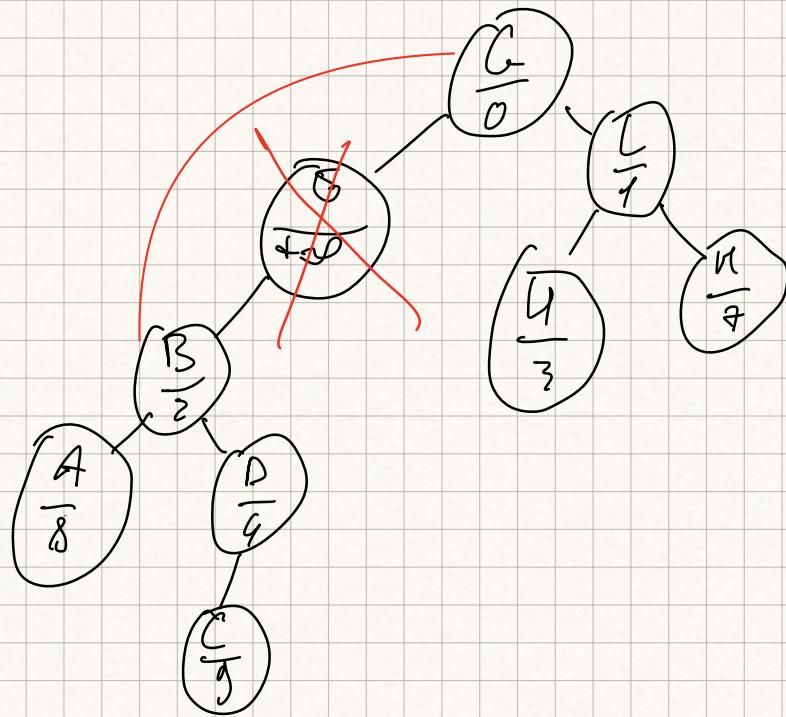
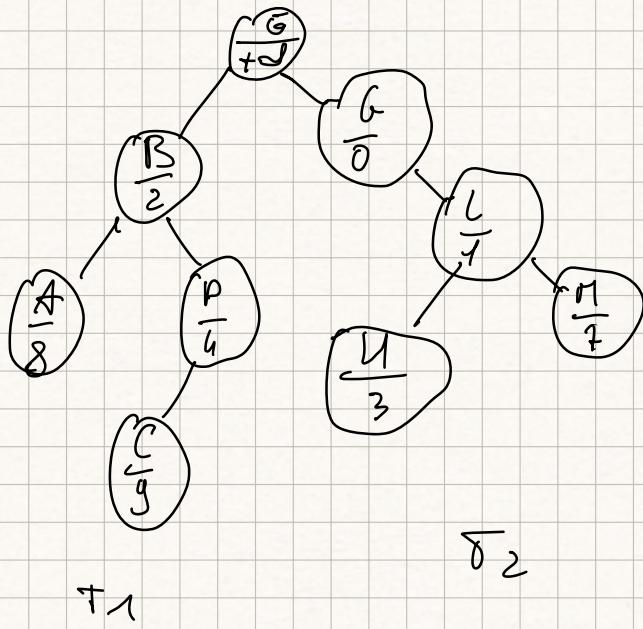
$$d'_1 = \boxed{1 \ 6 \ 15} \quad \Rightarrow \quad \{1, 2, 6\} \cap \{2, 3\} = 2$$

$$d_2 = \boxed{2 \ 3 \ 7} \quad \boxed{8 \ 18} \quad \Rightarrow \quad \{6, 9, 10\} \cap \{7, 8\} = \emptyset$$

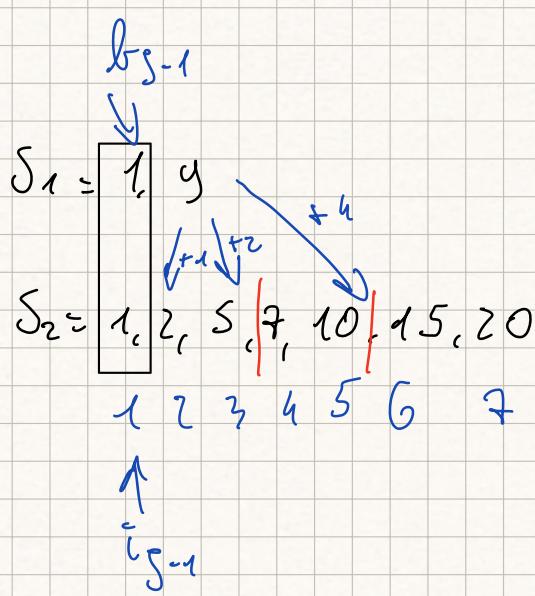
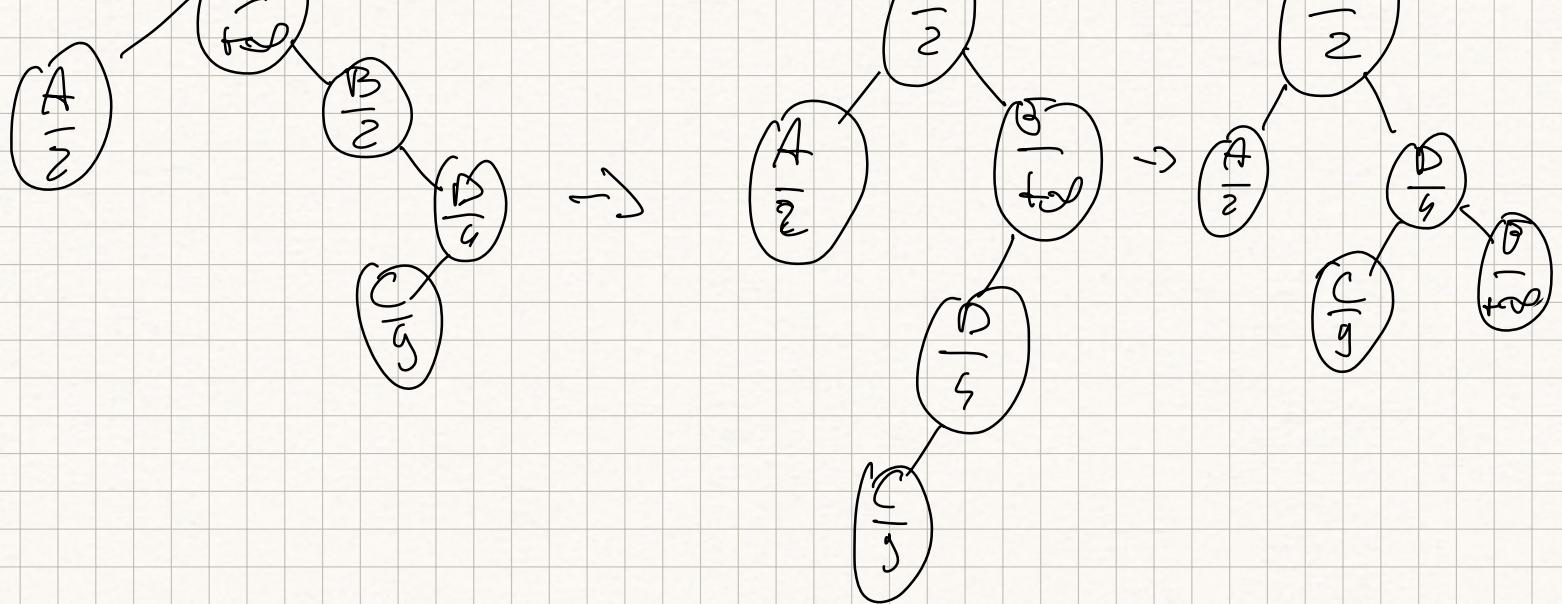
$$d'_2 = 2 \ 8 \quad \text{not used here} \quad \{15, 18, 20\} \cap \{18\} = 18$$

Merge d'_1 and d'_2 . In this case we do not skip any block

$$\begin{aligned} T_1 &= \{A, 8\}, \{B, 2\}, \{C, 9\}, \{D, 6\} \\ T_2 &= \{H, 3\}, \{I, 7\}, \{G, 0\}, \{L, 1\} \end{aligned} \quad \left. \right\} \text{ merge}$$



if I make continue



Exponential search

Mutual partitioning

$$S_1 = \boxed{1} g$$

$$S_2 = \boxed{1, 2, 5, 7, 10, 15, 20}$$

g

2, 5, 7, 10, 15, 20

$$k=3, \quad Q=2, \quad m=15$$

recycled items = $k(Q+1) - 1 = 8$

$$S = \{3, 1, 10, 9, 7, 6, 2, 50, 60, 14, 39, 21, 20, 5, 4\}$$

$$A = \{1, 3, 6, 2, 14, 39, 21, 4\}$$

$$= \boxed{1, 2, \textcircled{4}, 6, 9, \textcircled{14}} \quad 21, 39$$

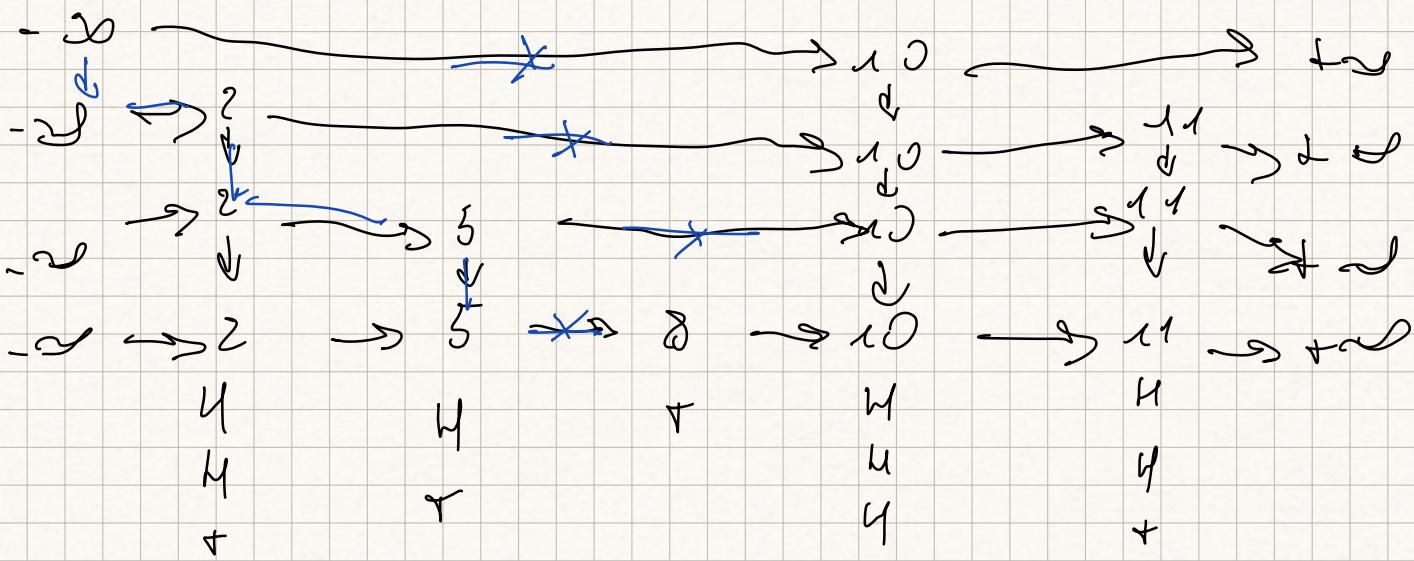
S_1 S_2

$$(-\infty, 4] \cup [9, 11] \cup (11, +\infty)$$

$$\begin{matrix} 3, 1, 2, 4 & 10, 9, 7 & 50, 60, 39, 21, 20 \\ 6, 14, 5 \end{matrix}$$

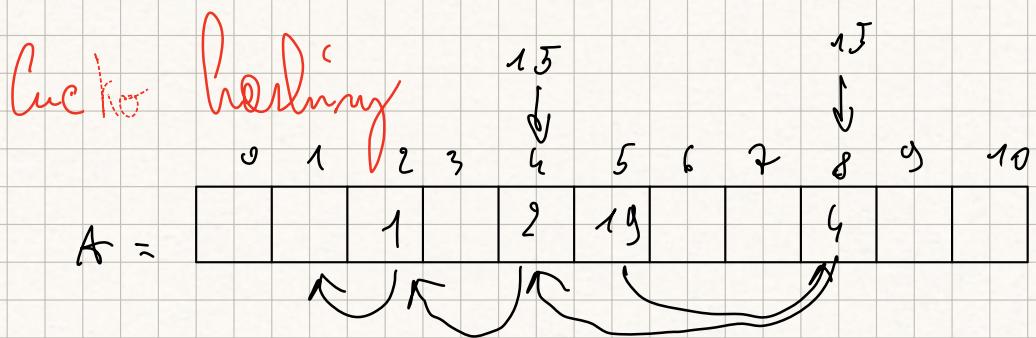
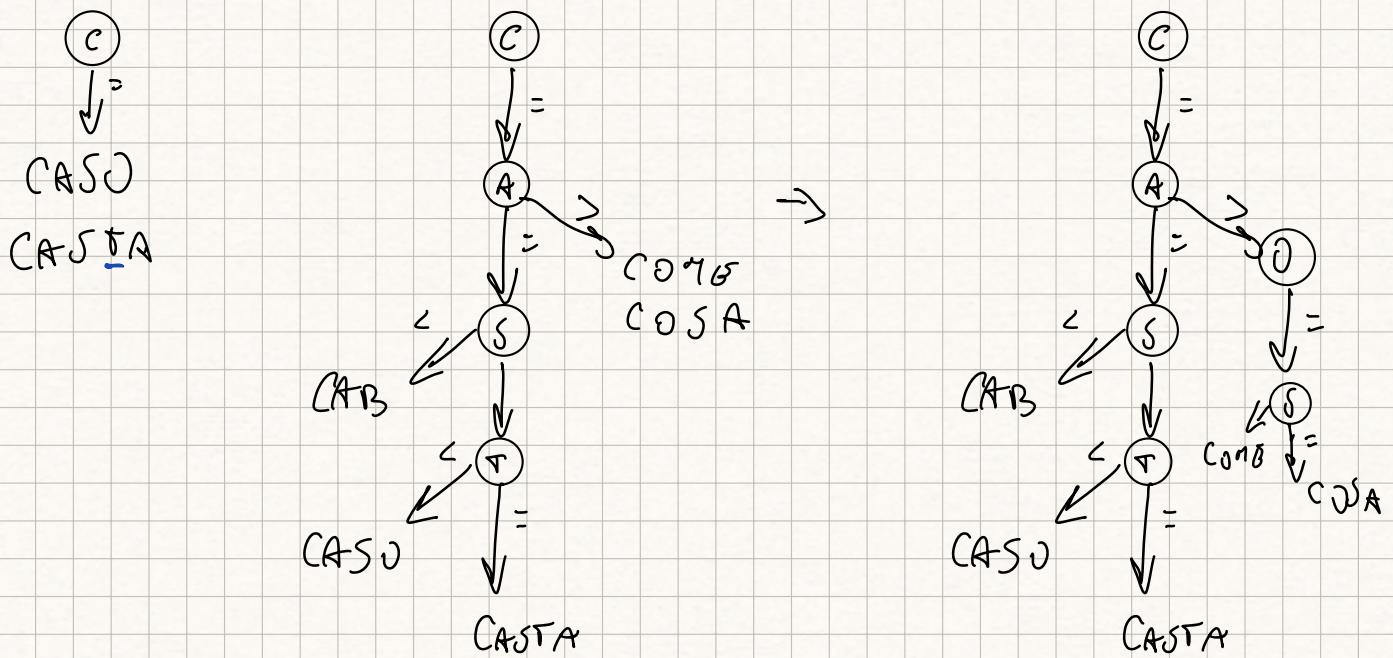
the. satisfied if bracket size is $\leq \frac{m}{k} = 20$

~~the. satisfied if bracket size is $\leq \frac{m}{k} = 20$~~



$$g = 6$$

$$S = \{ \text{CASTA}, \text{CASTA}, \text{CAB}, \text{CONG}, \text{COSA} \} \quad O(p + \log n) \text{ Time}$$

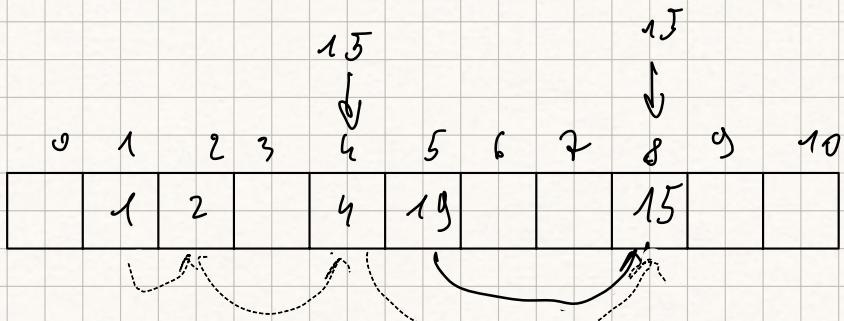


$$m = 11$$

$$S = \{ 1, 4, 2, 19, 15 \}$$

$$h_1(x) = 2x \bmod 11$$

$$h_2(x) = x \bmod 11$$



Given the two sorted lists $L_1 = (1, 2, 4, 6, 9, 10, 15, 18, 20)$ and $L_2 = (2, 3, 7, 8, 18)$ compute their intersection using the

1) Mutual partitioning strategy

2) Two-level storage approach, with block size $b=3$ for the list L_1

1)

$$L_1 = \underbrace{1, 2, 4, 6}_{L_{11}}, \underbrace{9, 10, 15, 18, 20}_{L_{12}}$$

$$L_2 = \underbrace{2, 3}_{B_1}, \underbrace{7, 8}_{B_2}, \underbrace{18}_{B_3}$$

$$L_{11} \cap L_{21} = \{1, 2, 4, 6\} \quad L_{11}' = \{3\}$$

$$L_{12} \cap L_{22} = \{9, 10, 15, 18, 20\} \quad L_{12}' = \{18\}$$

$$L_{11}' \cap L_{21}' = \{3\}$$

$$L_{12}' \cap L_{22}' = \{18\}$$

$$L_1 = \underbrace{1, 2, 4, 6}_{B_1}, \underbrace{9, 10, 15}_{B_2}, \underbrace{18, 20}_{B_3}$$

$$L_2 = \underbrace{2, 3}_{B_1}, \underbrace{7, 8}_{B_2}, \underbrace{18}_{B_3}$$

$$L_1' = 1, 6, 15$$

Merge L_2 and L_1'

$$B_1 \cap L_{11} = 2$$

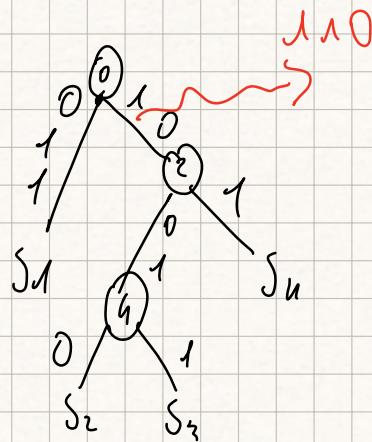
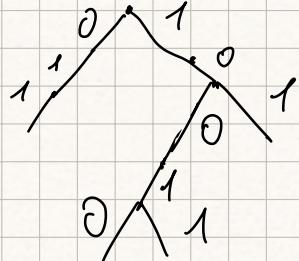
$$B_2 \cap L_{12} = \emptyset$$

$$B_3 \cap L_{13} = 18$$

Given the binary strings $S = \{011, 10010, 10011, 101\}$.

- Build the Patricia Trie for S

- Show how to search for the lexicographic position of the string $P=110$ among the strings of S .



Given the sequence of 6 items $S = (a, b, c, d, e, f)$ use the random sampling algorithm with known sequence length $n=6$, to compute the $m=2$ extracted items given the random values $p = (1/2, 1/10, 3/4, 3/4, 0, 1)$

$S = a \overset{1}{b} \overset{2}{c} \overset{3}{d} \overset{4}{e} \overset{5}{f} \overset{6}{}$ $m = 2$

$$1/2 \quad 1/10 \quad 3/4 \quad 3/4 \quad 0 \quad 1$$

$$P \leftarrow \underbrace{m-j}_{m-j+1}$$

$$j=0 \quad j=1 \quad \frac{1}{2} < \frac{2-0}{6-1+1} = \frac{1}{3} \times$$

$$j=1 \quad j=5 \quad 0 < \frac{5-1}{6-5+1} = \frac{1}{2} \checkmark$$

$$j=0 \quad j=2 \quad \frac{1}{10} < \frac{2-0}{6-2+1} = \frac{2}{5} \checkmark$$

$$j=1 \quad j=3 \quad \frac{3}{4} < \frac{3-1}{6-3+1} = \frac{1}{3} \times$$

$$j=1 \quad j=4 \quad \frac{3}{4} < \frac{4-1}{6-4+1} = \frac{1}{3} \times$$

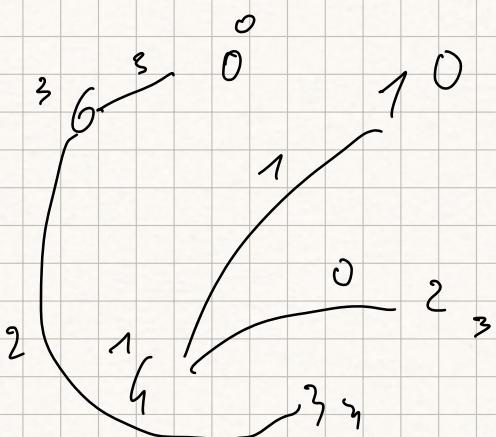
	h	h_1	h_2
AA	0	4	2
BB	1	1	4
BD	2	3	6
CD	3	6	0

$$h_1 = 3 \text{ rank}(x) + \text{rank}(y) \bmod 7 \quad \checkmark$$

$$h_2 = \text{rank}(x) + 2\text{rank}(y) \bmod 7 \quad \checkmark$$

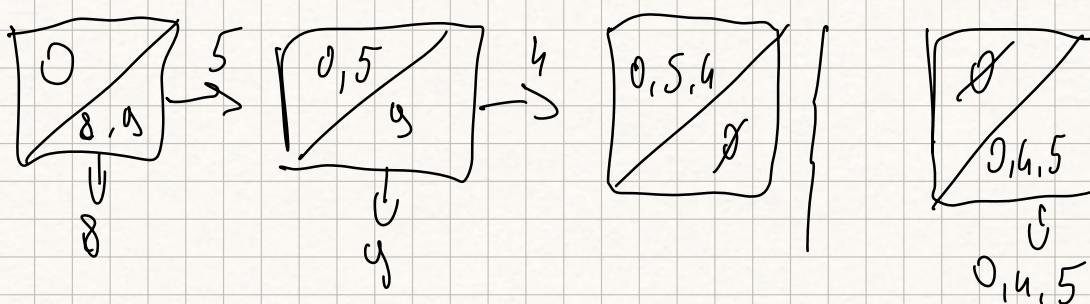
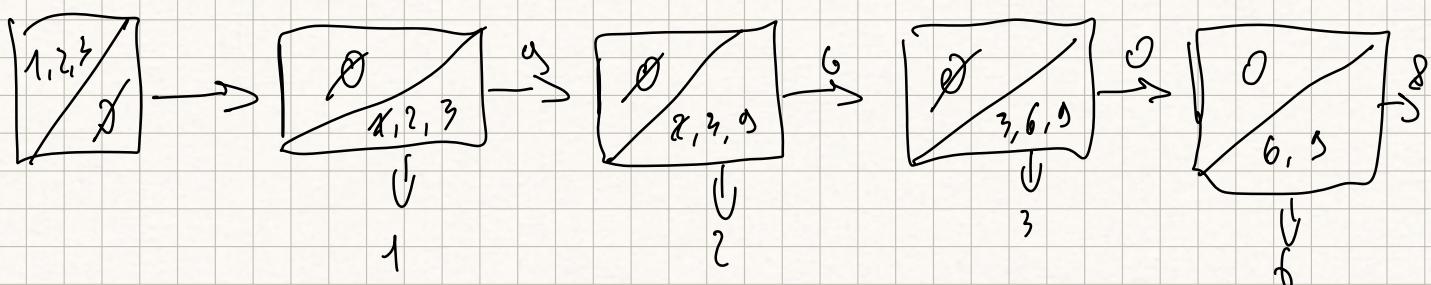
$\text{rank}(1, 2, 3, 6)$

$$h = g(h_1) + g(h_2) \bmod 4$$

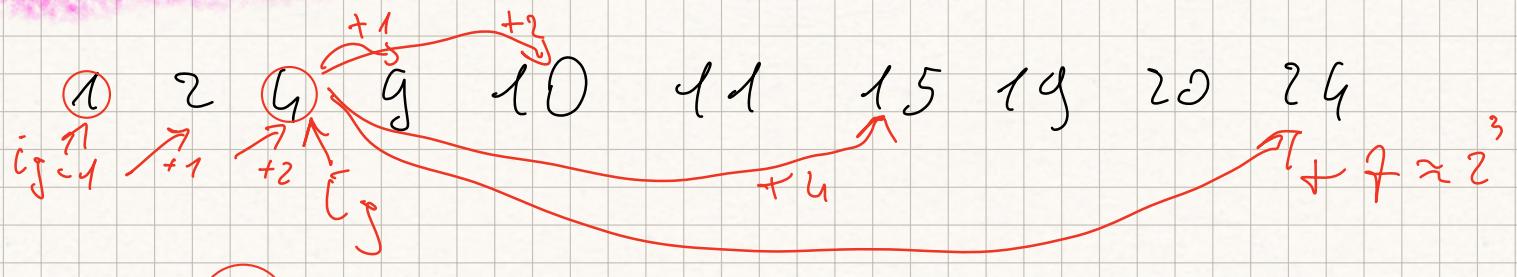
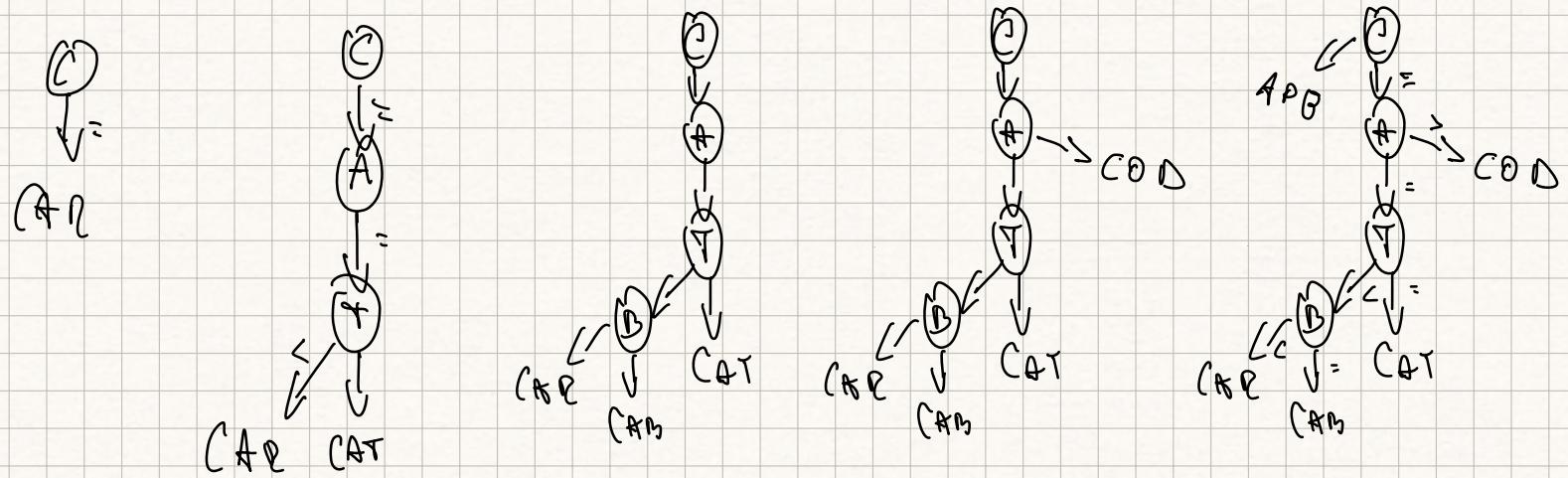


$$g = 0 \ 0 \ 3 \ 3 \ 1 \ 0 \ 3$$

Snowflake : 1, 2, 3, 9, 6, 0, 8, 5, 4 $B = 3$



Build a Ternary Search Tree by inserting the following strings according to the given sequence (CAR, CAT, CAB, COD, APE)



focus use binary search tree 9, 14 e next tree 3

height and binary search tree 15, 24 e tree 20

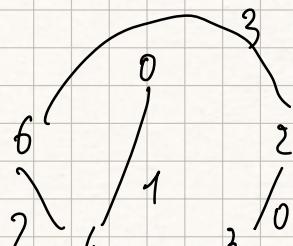
	h	h_{11}	h_{12}
AA	0	2	3
AC	1	4	0
BB	2	4	6
CC	3	6	2

$$h_{11} = x + y \bmod 7$$

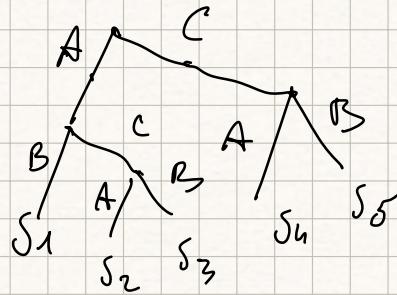
$$h_{12} = x + 2y \bmod 7$$

$$h = g(h_{11}) + g(h_{12}) \bmod 4$$

$$g = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ 0 & 0 & 2 & 2 & 4 & 0 & 1 \end{matrix}$$



AB, ACA, ACB, CA, CB



$$A = \underbrace{3 \ 5 \ 6}_{A_1} \ \underbrace{8 \ 10 \ 13}_{A_2} \ \underbrace{14 \ 16 \ 17}_{A_3} \quad B = 3 \ 7 \ 8 \ 14 \ 15$$

$B = 3$

$B = \underbrace{5 \ 7 \ 8}_{A_1 \ A_2} \ \underbrace{14 \ 15}_{A_3}$

merge $\underline{3} \ 5 \ 7 \ 8 \ \underline{8} \ 14 \ 15$

$A' = 3 \ 8 \ 14$

Q5: CASTRO, ABBA, NOM, CAMB, ASTRA, ASSO

$$P = ABBA$$

$$P = ABBA$$

$R_L = ABBA, ASTRA, ASSO$

$$\begin{array}{c} i=1 \\ R_L = \left\{ \begin{array}{l} ABBA \\ ASTRA \\ ASSO \end{array} \right. \end{array} \xrightarrow{i=2} \begin{array}{c} Q_L = ABBA \\ R_L = ASTRA \\ ASSO \end{array}$$

$$P = ASTRA$$

$$P_L = ASSO$$

$R_L = CASTRO, CAMBL$

$$\xrightarrow{i=2}$$

$$P = CASTRO$$

$R_L = \boxed{NOM}$

$$\begin{array}{c} i=3 \\ R_L = \left\{ \begin{array}{l} CASTRO \\ CAMBL \end{array} \right. \end{array} \xrightarrow{i=3} \begin{array}{c} R_L = \boxed{CAMBL} \\ P_L = \boxed{CASTRO} \end{array}$$

$S = 1, 6, 15, 18, 21, 24, 30$

↓

$S = 1, 5, 8, 2, 3, 3, 6$

gap encoding

f - code

f - code

Rice - code

$k=3$

(1, 3) - dense code

f - code (s') = 1 00101 0001001 011 011 011 00110

f - code = f code (e) bin(x) = 11 01101 0010001001 01011 01011 00011 01110

Rice - code :

Unary Binary,

1 100

R - c (5) $\Rightarrow q = 0$ $q = 4$

R - c (3) $\Rightarrow q = 1$ $q = 2$

R - c (3) $\Rightarrow q = 2$ $q = 2$

R - c (6) $\Rightarrow q = 0$ $q = 5^{101}$

$r = x-1 - 2^k \cdot q$

Dense - code :

	config			
1	00	5		
2	01 00	c 5		
3	10 00	c 5		
4	11 00	c 5		
5	01 01 00	c c 5		
6	01 10 00	$3 \times 3 \times 1$		
7	01 11 00	$= 3$		
8	10 01 00	$c c c 5 = 27 > 25$		
9	10 10 00	$3 \times 3 \times 3 \times 1$		

$$1+3=4=2^2 \leftarrow b$$

$$\begin{array}{r} 00 5 \\ 01 \\ 10 \\ 11 \end{array}$$

= 3

$$c c c 5 = 27 > 25$$

Blis - Faw:

$$H = 0110101$$

$$L = \overbrace{1}^0 \overbrace{0}^1 \overbrace{1}^2 \overbrace{0}^3 \overbrace{0}^4 \overbrace{1}^5 \overbrace{1}^6 \overbrace{0}^7$$

$$\mu = \tau = \#\text{1 in } H$$

$$l = d = \frac{|L|}{\mu}$$

$$k = 3 \quad \#\text{0 in } H$$

$$b = 4$$

	H	L
0	000	0
3	001	1
5	010	1
8	100	0
9	100	1
10	101	0
11	101	1

1, 3, 5, 9, 16, 23, 27, 28, 31, 40

$$\mu = 11$$

$$d = 61$$

$$l = \lceil \log_2 \frac{m}{\mu} \rceil = 2$$

$$b = \lceil \log_2 m \rceil = 6$$

$$h = 6-2 = 4$$

	H	L
1	00000	01
3	00000	11
5	000100	
9	000000	1
16	010000	
23	010111	
27	011011	
28	011100	
31	011111	
40	101000	

$$L = 01100001010011100100$$

$$H = 1011010010101010110000000$$

$$\text{Access}(u) = 000101 \quad \text{Select}_x(u) = 5-4 = 1$$

$$\text{Get}_R(8) = ?$$

$$\text{Get}_L(2) = ?$$

$$\text{Access}(7) = 010111$$

$$\text{Select}_x(7) = 12-7 = 5$$

$\text{GCD}(8) = 0010\overset{1}{00}$ Select₀(2) + 1 = 6 + 1 = 7 select the right bracket

$$7 - 2 = 5$$

$$\text{Acan}(5) = 9$$

$\text{GCD}(32) = \underbrace{1000}_{8}00$ $P = \text{select}_0(8) + 1 = 18 + 1 = 19$

$$P - 8 = 11$$

$$\text{Acan}(11) = 40$$

Interpolating code

$$S = 11, 19, 16, 10, 20, 21, 22$$

$$m = \lceil \frac{l+r}{2} \rceil = 4 \quad S_m = 19$$

l = 1
r = 7
low = 11
new = 22

$$Q = \text{low} + m - l = 11 + 4 - 1 = 19$$

$$R = \text{high} + m - r = 22 + 4 - 7 = 29$$

encode $S_m - Q = 5$ in $\lceil \log_2 R - Q + 1 \rceil = 3$ (01)₂

$$l = 5, r = 7, \text{low} = 20, \text{hi} = 22 \Rightarrow m = 6, S_m = 21$$

right

$$\left. \begin{array}{l} Q = 20 + 6 - 5 = 21 \\ R = 22 + 6 - 7 = 21 \end{array} \right\} \text{encode } 0 \text{ in } 0 \text{ bit}$$

$\lceil \log_2 17 \rceil$

left

$$l = 1, r = 3, low = 11, high = 18$$

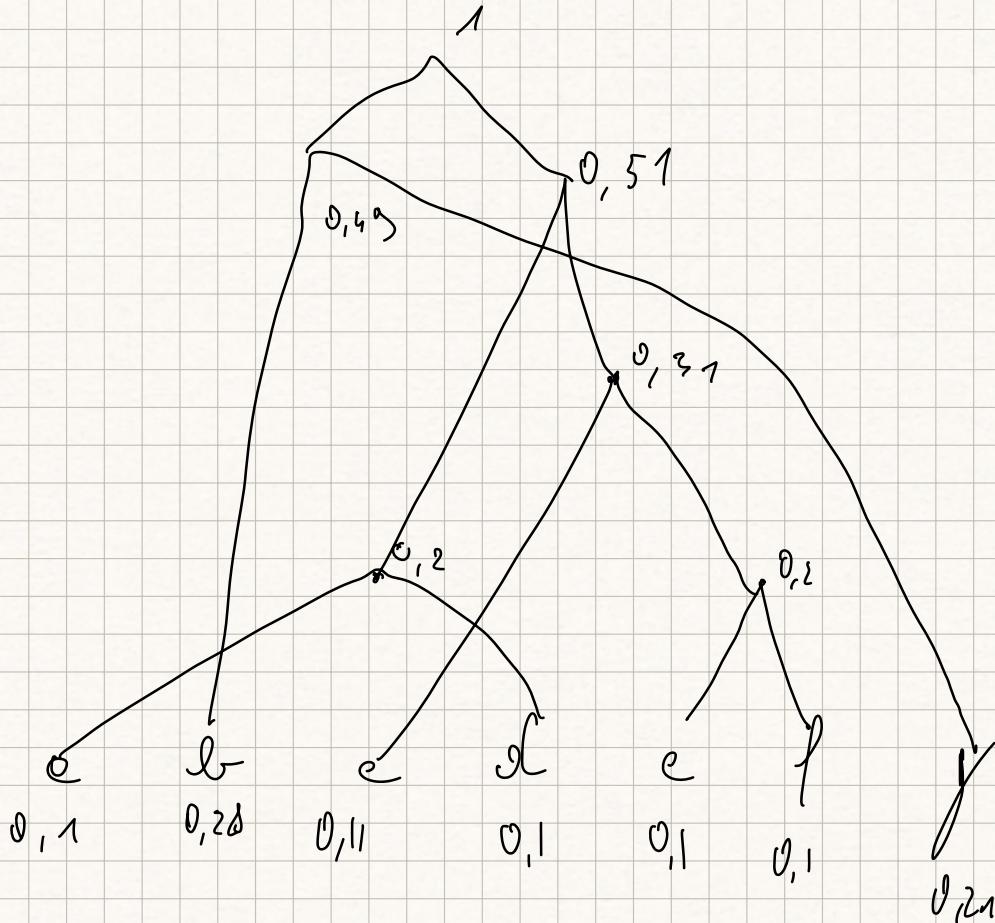
$$m = 2 \quad S_m = 14$$

$$Q = 11 + 2 - 1 = 12$$

$$B = 18 + 2 - 3 = 17$$

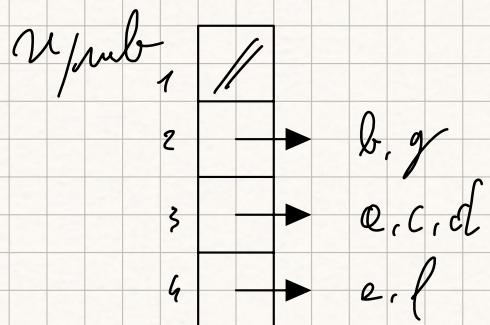
encode 2 in $\lceil \log_2 17 - 1 \rceil + 1 = 3$ $(0|0)_2$

- Huffman tree
- Canonical Huffman tree
- Decompress 11001 using Canonical H.T.



$$\min[1, 4] = \begin{bmatrix} 0, 1, 3, 2 \\ 1 \quad 2 \quad 3 \quad 4 \end{bmatrix}$$

$$FC[4] = 0$$



$$FC[3] = \frac{FC[4] + \min[4]}{2} = 1$$

$$FC[2] = \frac{FC[3] + \min[3]}{2} = 2$$

$$FC[1] = \frac{FC[2] + \min[2]}{2} = 2$$

$$FC[1, 4] = [2, 2, 1, 0]$$

Decompressor:

$$v = \text{next_bit}();$$

$$l = 1;$$

$$\text{while } (v < FC[l])$$

$$v = 2 \cdot v + \text{next_bit}();$$

$$l++$$

$$\text{return symbol}[l, v - FC[l]]$$

11001

$$v = 1 \quad l = 1 \quad v = 1 < FC[1] = 2$$

$$v = 3 \quad l = 2 \quad v = 3 > FC[2] = 2$$

$$\text{symbol}[2, 1] = g$$

$$v = 0 \quad l = 1 \quad v = 0 < FC[1] = 2$$

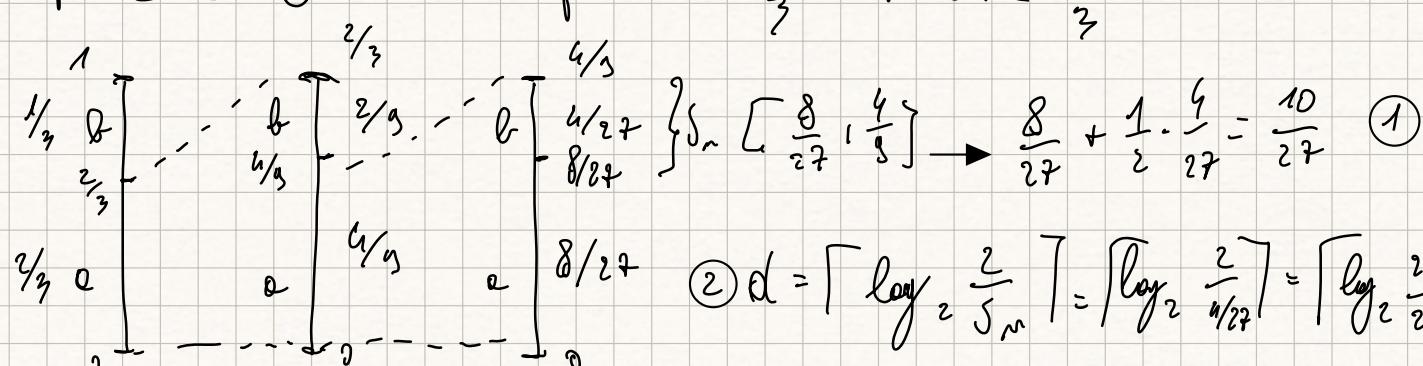
$$v = 0 \quad l = 2 \quad v = 0 < FC[2] = 2$$

$$v = 2 \cdot 0 + 1 = 1 \quad l = 3 \quad v = 1 < FC[3] = 1$$

$$\text{symbol}[1, 1-1] = a$$

$$T = a \quad e \quad b$$

$$P(a) = \frac{2}{3} \quad P(b) = \frac{1}{3}$$



$$\sum_{n=1}^{\infty} \left[\frac{8}{27}, \frac{9}{27} \right] \rightarrow \frac{8}{27} + \frac{1}{2} \cdot \frac{9}{27} = \frac{10}{27} \quad (1)$$

$$(2) d = \lceil \log_2 \frac{2}{\frac{10}{27}} \rceil = \lceil \log_2 \frac{2}{\frac{10}{27}} \rceil = \lceil \log_2 \frac{27}{10} \rceil = 4$$

$$S_0 = 1$$

$$S_1 = \frac{1}{3}, \quad S_2 = \frac{4}{9}$$

$$Q_0 = 0$$

$$Q_1 = 0, \quad Q_2 = 0$$

③ Converter $\left(\frac{10}{27}, 4\right) = 0101$

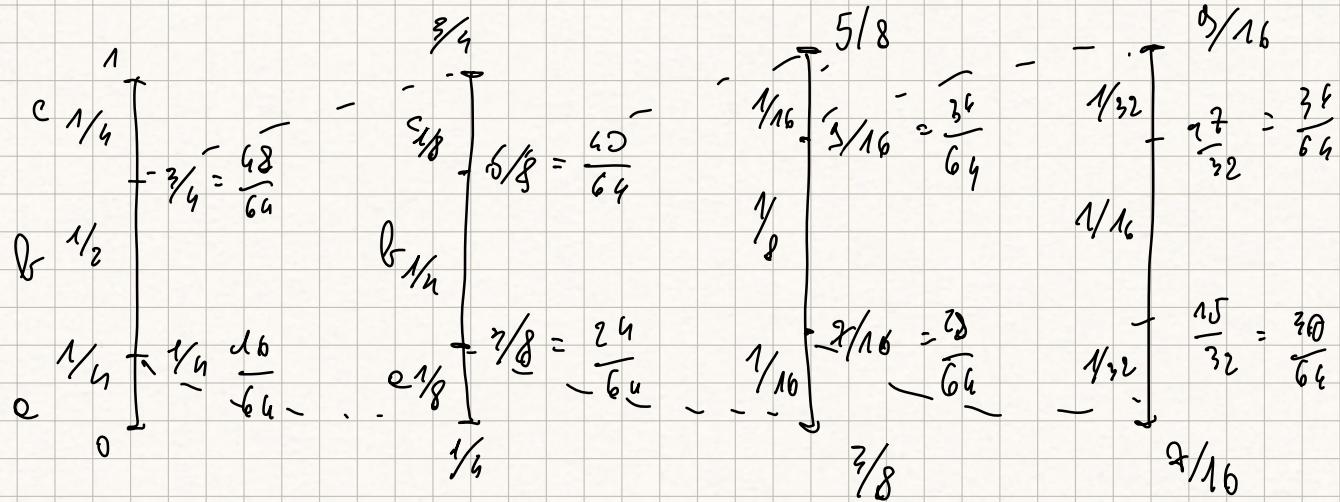
$$\frac{10}{27} \rightarrow \frac{20}{27} < 1 \rightarrow 0$$

$$\frac{20}{27} \rightarrow \frac{40}{27} > 1 \rightarrow 1 \rightarrow \frac{13}{27}$$

$$\frac{13}{27} \rightarrow \frac{26}{27} < 1 \rightarrow 0$$

$\langle 011111, 4 \rangle \quad P(a) = P(c) = \frac{1}{4}, P(b) = \frac{1}{2} \rangle$

$$\frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \frac{1}{32} + \frac{1}{64} = \frac{16+8+4+2+1}{64} = \frac{31}{64}$$



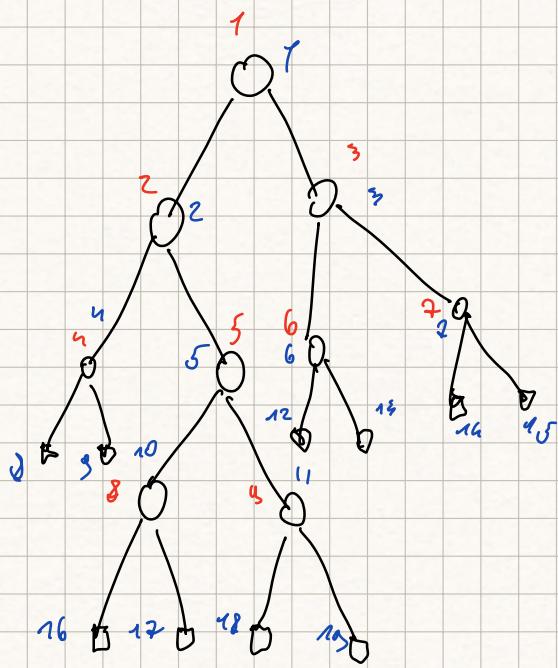
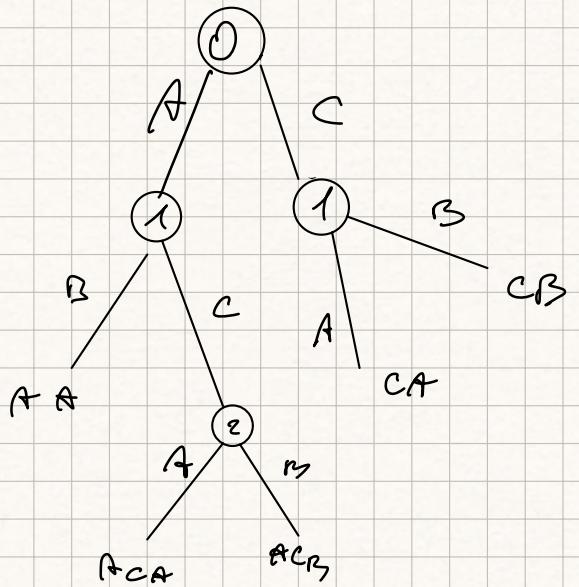
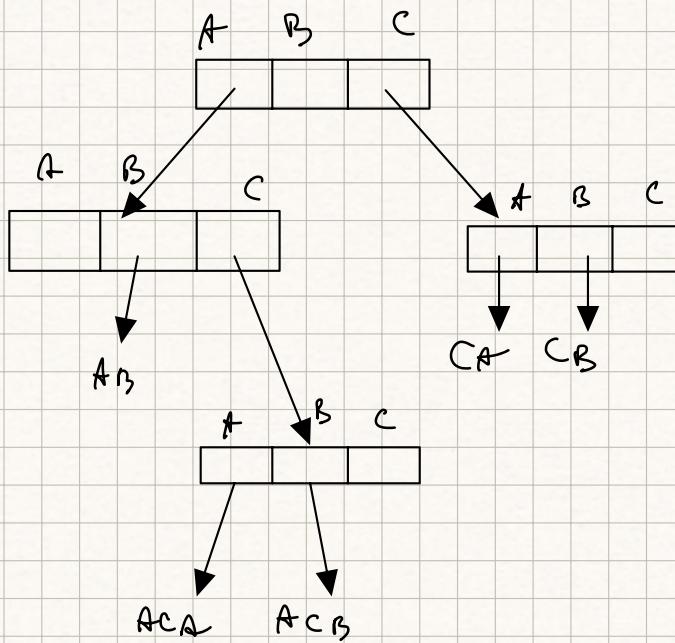
$$S_1 = 1/2$$

$$Q_1 = 1/4$$

$$S_2 = 4/9$$

$$Q_2 = 3/8$$

$D = \{AB, AC\bar{A}, \bar{A}CB, CA, C\bar{B}\}$ Succinctly encode the tree



0 dummy nodes

1 real nodes

$$B = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 \\ 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 \end{matrix}$$

$$\begin{matrix} 16 & 17 & 18 & 19 \\ 0 & 0 & 0 & 0 \end{matrix}$$

Pseudocode: return the length of the leftmost path

$\text{Rank}_L(zx)$

$\Rightarrow [zx] + 0$

$left = 0$

$x = 1$

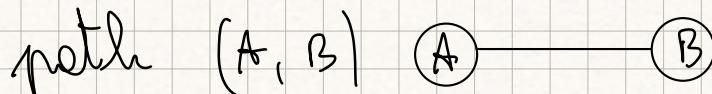
while ($B[2x] \neq 0$) do

$left++;$

$x = rank_1(2x)$

return $left$

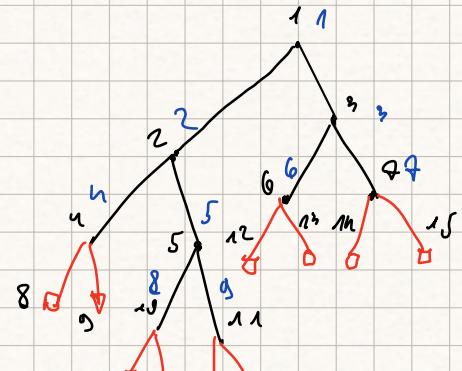
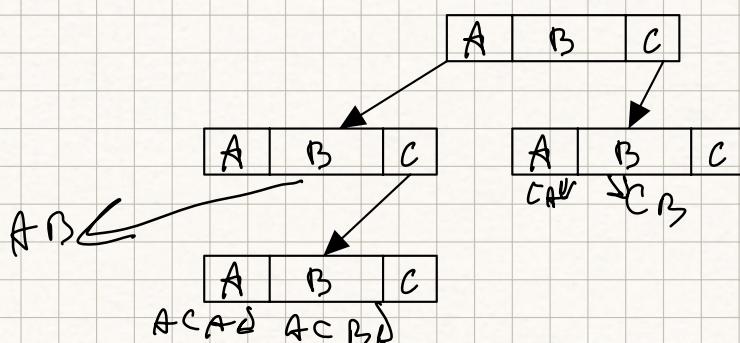
write a pseudo-code that check if the tree contains a path (A, B)



Question #2 [scores 4+3+3] Given the ordered set of strings

$$S = \{AB, ACA, ACB, CA, CB\}$$

- Build the UNcompacted trie T for S by assuming an alphabet of 3 characters ($S = \{A, B, C\}$) and branching implemented via arrays.
- Show how to succinctly encode the structure of T in a binary array B, by assuming that pointers to strings are leaves of the tree T, and branching nodes are the internal nodes of T;
- Write a pseudo-code that, given a binary array succinctly encoding the structure of a binary tree (with its corresponding rank/select data structures), establishes the length of its left-only path (or, equivalently, the depth of its leftmost NULL pointer).



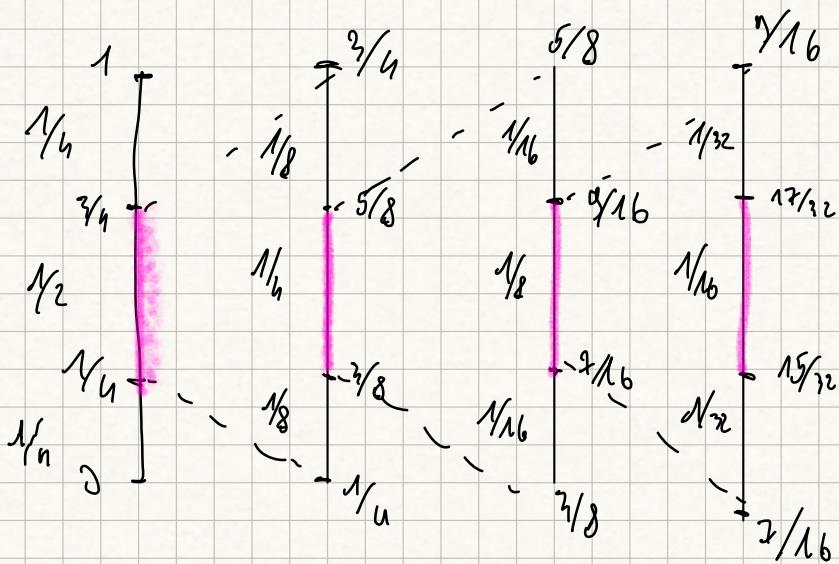
16 17 18 19

$B = 11111110011000000000$

Decode the compressed sequence $<4, 011110>$ produced by arithmetic code, by assuming probabilities $P[a]=P[c]=1/4$ and $P[b]=1/2$.

$$0 \ 1 \ 1 \ 1 \ 1 \ 0 \Rightarrow \frac{1+2+4+8}{32} = \frac{15}{32} \approx \frac{1}{2}$$

$$\frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \quad \frac{1}{16} \quad \frac{1}{32} \quad \frac{1}{64}$$



Show the first 8 codewords of the (s,c) -code with $s=3$ and $c=1$, hence $s+c = 4$ (briefly explain your calculations).

$$S = 00 \quad C = 11$$

01

10

Codeword (3,1)

1	00
2	01
3	10
4	1100
5	1101
6	1110
7	111100
8	111101

$$S+C=4=2^B \Rightarrow B=2$$

Decompress the 8th integer encoded via Elias-Fano in the two arrays:

$L = 01|11|00010100|11|1100|1100$ and $H = 110\ 110\ 10\ 0\ 10\ 10\ 10\ 110\ 0\ 0\ 10\ 0\ 0\ 0\ 0\ 0$ where the original encoding of the integers is in 6 bits. (hint: derive first the number of keys, and then the length of the low and high part)

$$\text{Access}(8) : \text{ take the } 8\text{-th group in } L = \left. \begin{array}{c} 11 \\ 0110 \\ 11 \end{array} \right\} \frac{h}{0110} \frac{l}{11} = 27$$

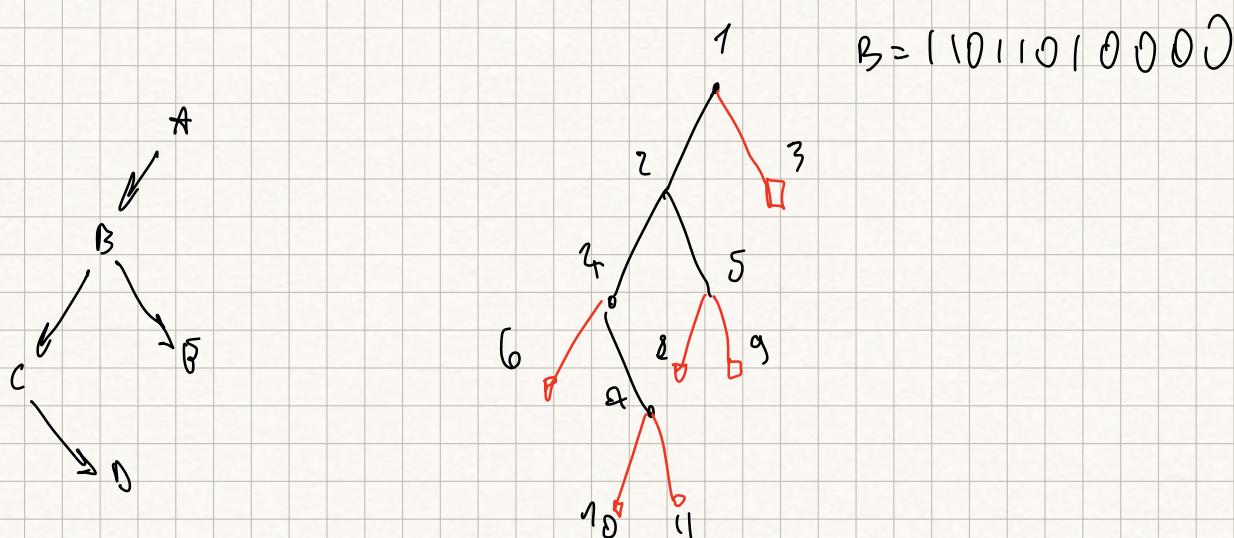
$$\text{Select}_1(8) = 19 - 8 : (6)_2 = 110$$

$$h = 6$$

$$l = L / \# 1 \text{ in } H = 2$$

$$k = 4$$

[rank 4]. Show the binary succinct encoding of the tree $T = \{ \text{aab} \text{ (left child)}; \text{bac} \text{ (left child)}; \text{bae} \text{ (right child)}; \text{cad} \text{ (right child)} \}$ of root "a".

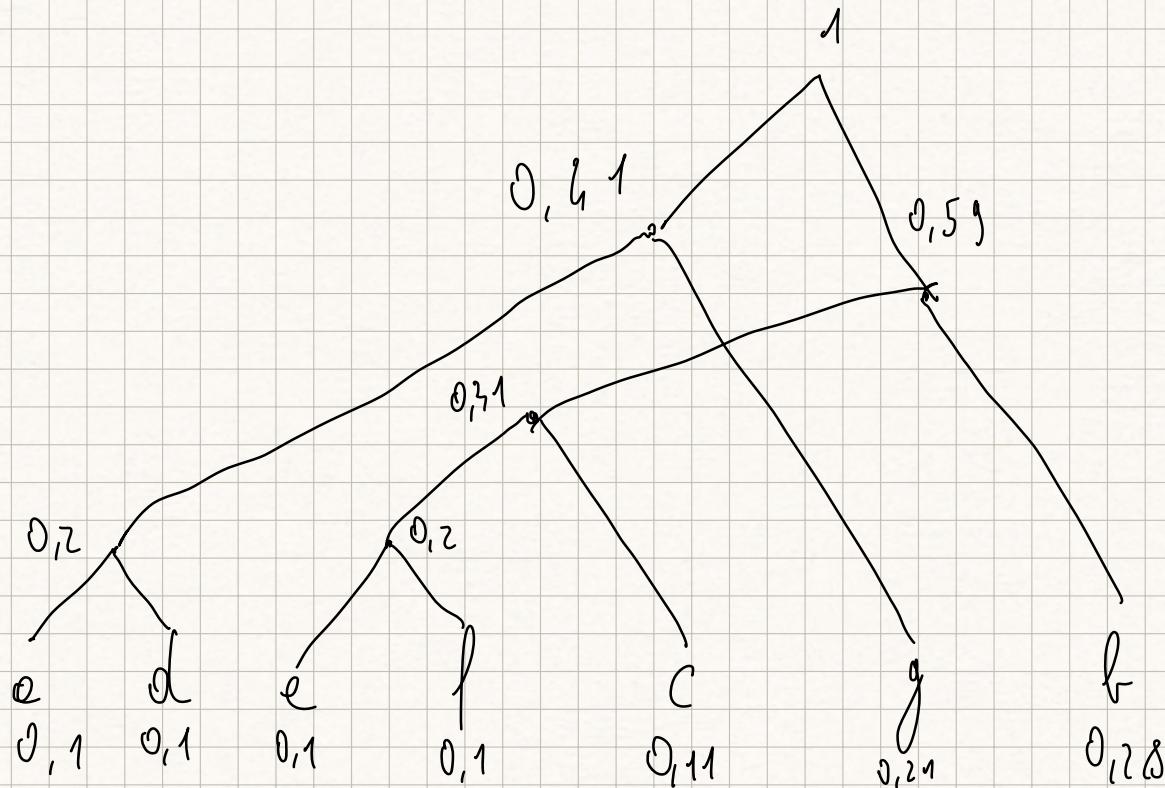


Given the symbols {a,b,c,d,e,f,g} occurring in a text with frequencies $f(a) = f(d) = f(e) = f(f) = 0.1$, $f(b) = 0.28$, $f(c) = 0.11$, $f(g) = 0.21$.

- Compute FC[] and SYMB[] tables of the Canonical Huffman code

• Decode the first 2 symbols of the compressed sequence:

11001...



$$\text{num}(1, h) = [0, 2, 3, 1]$$

$$FC[i] = \frac{FC[i+1] + \text{num}[i+1]}{2}$$

1 1 | 0 0 1

$v = \text{first_bit}()$ $l = 1$

while $v < FC[e]$

$v = 2 \cdot v + \text{next_bit}()$

$l++$

return symb($l, v - FC[e]$)

SYMB	FC
1	
2	
3	
4	

$$v=1 \quad l=1 \quad 1 < 2$$

$$v=2 \cdot 1 + 1 = 3 \quad l=2 \quad 3 > 2$$

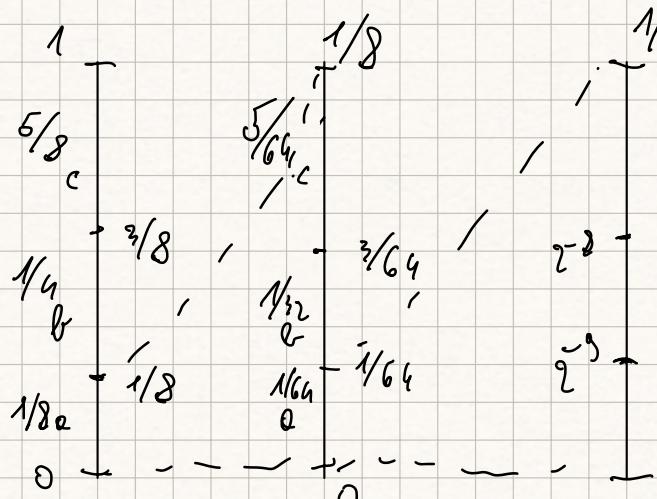
$$\text{symb}[2, v - FC[1]] = g$$

$$v=0 \quad l=1 \quad 0 < 2 \quad \text{symb}[3, 0] = 0$$

$$v=2 \cdot 0 + 0 \quad l=2 \quad 0 < 2$$

$$v=2 \cdot 0 + 1 \quad l=3 \quad 1 \geq 1$$

Given: $p(a) = 1/8$, $p(b) = 1/4$, $p(c) = 5/8$. Specify which is the length in bits of the text $T = aabbba$, if it is compressed via Arithmetic coding. (Hint: work with negative powers of two.)

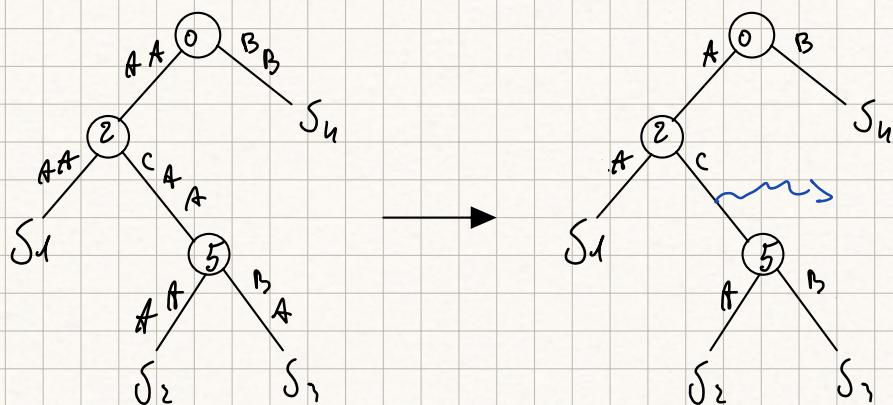


$$d = \left\lceil \log_2 \frac{2}{S_m} \right\rceil$$

$$S_m = \left(\frac{1}{8}\right)^4 \cdot \left(\frac{1}{4}\right)^2 = \frac{1}{2^{12}} \cdot \frac{1}{2^4} = \frac{1}{2^{16}}$$

$$d = \left\lceil \log_2 \frac{2}{2^{-16}} \right\rceil = \left\lceil \log_2 2^{16} \right\rceil = 16 \text{ bits}$$

Given the binary strings $S = \{\text{aaaaa}, \text{aacaabaa}, \text{aacaaba}, \text{bb}\}$. Build the Patricia Tree for S and show how to search for the lexicographic position of the string $P = \text{aacbba}$ among the strings of the set S .



$P = A A C B B A$ search S_2 , where $\text{LCP}[P, S_2] = \frac{A A C}{3}$, quindi corri

Patricia tree S_3 e S_n

Given the integer sequence $S = (1, 2, 3, 4, 6, 8, 9)$, show how Interpolative Coding compresses the "first three" integers according to its algorithm.

$$(l=1, r=7, \text{low}=1, \text{high}=9)$$

$$m = \left\lfloor \frac{l+r}{2} \right\rfloor = 4 \quad S_m = 4$$

$$a = \text{low} + (m-l) = 1 + 3 = 4$$

$$b = \text{high} - (r-m) = 9 - 3 = 6$$

$$S_m - a = 4 - 4 = 0 \quad (00)_2$$

$$\lceil \log_2(b-a+1) \rceil = 2$$

$$(l=1, r=3, \text{low}=1, \text{high}=3)$$

$$m = \left\lfloor \frac{l+r}{2} \right\rfloor = 2 \quad S_m = 2$$

$$a = \text{low} + m-l = 1 + 1 = 2$$

$$b = \text{high} - (r-m) = 3 - 1 = 2$$

$$(l=5, r=7, \text{low}=5, \text{high}=9)$$

$$m = \left\lfloor \frac{l+r}{2} \right\rfloor = 6 \quad S_m = 8$$

$$a = \text{low} + m-l = 5 + 6 - 5 = 6$$

$$b = \text{high} - (r-m) = 9 - 1 = 8$$

$$S_m - a = 0 \quad \lceil \log_2(b-a+1) \rceil = 1$$

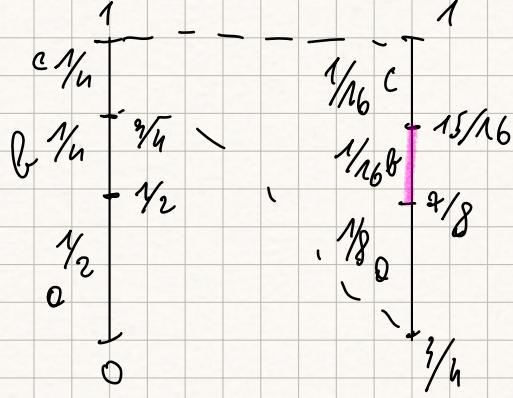
$$S_m - a = 2 \quad \lceil \log_2(b-a+1) \rceil = 2 \quad (10)_2$$

more about 5 symbols

Let us given the probabilities: $p(a) = 1/2$, $p(b)=p(c)=1/4$.

Decompress the first 2 symbols of the Arithmetic coded bit sequence: 111.

$$\begin{array}{ccc} 1 & 1 & 1 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{8} \end{array} = \frac{7}{8}$$



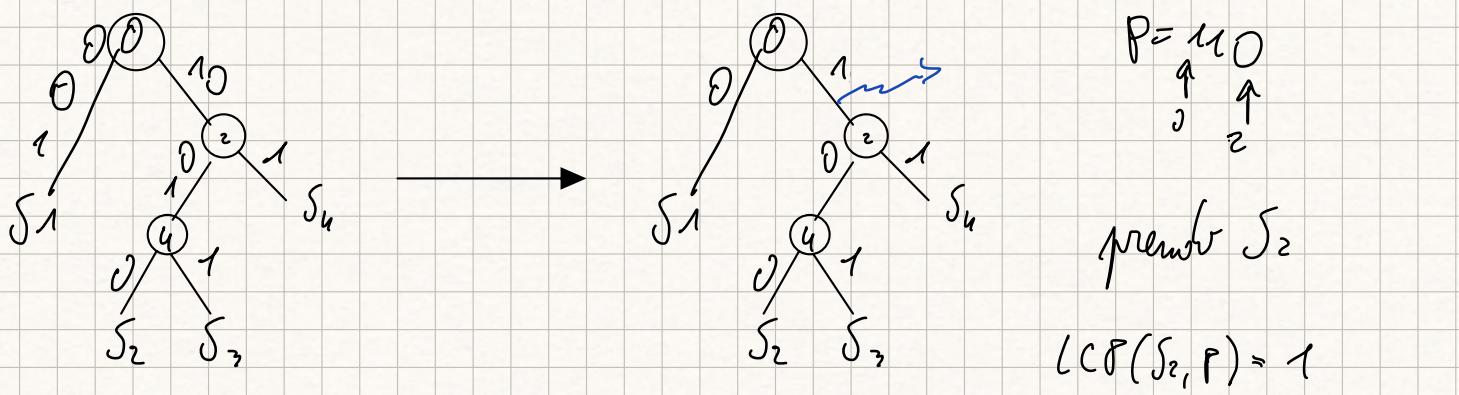
$$S_0 = 1 \quad S_1 = \frac{1}{4}$$

$$\ell_0 = 0 \quad \ell_1 = \frac{3}{4}$$

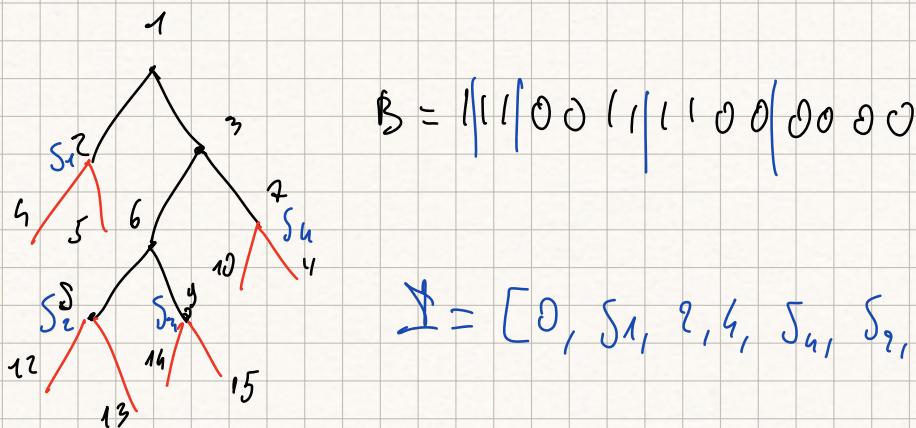
c b

Given the binary strings $S = \{001, 10010, 10011, 101\}$.

- Build the Patricia Trie for S
- Show how to search for the lexicographic position of the string $P=110$ among the strings of the set S .
- Propose a succinct encoding of the Patricia Trie of S that allows navigation in constant time per traversed edge.
- Simulate the downward search for $P=110$ in this succinct encoding.



herw 0 obtre di S_4



] Let us given the set of strings

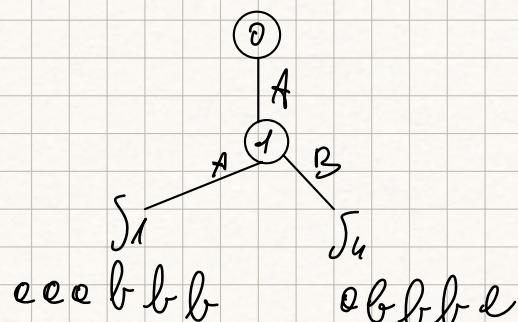
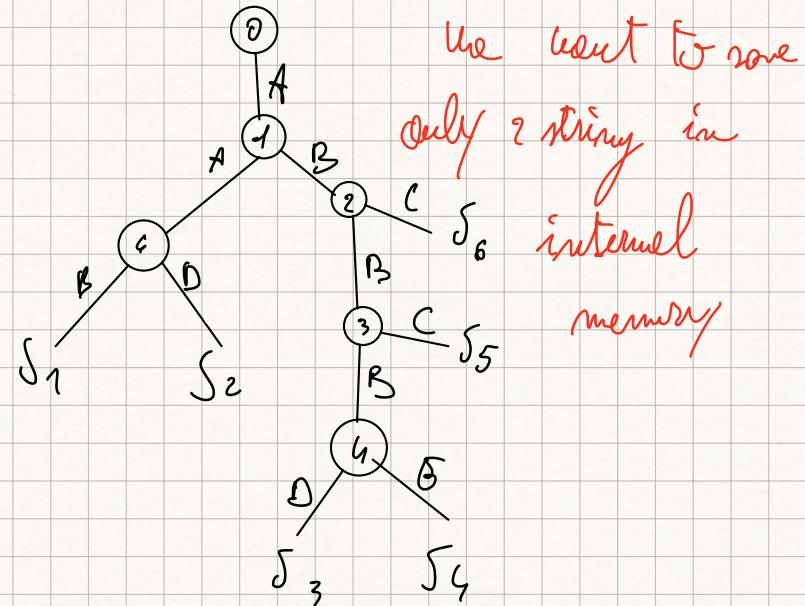
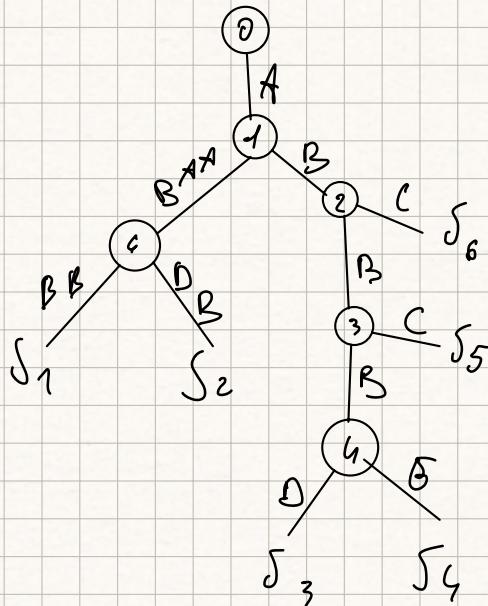
$S = \{ \text{baci}, \text{baco}, \text{boss}, \text{buono}, \text{buu}, \text{costi}, \text{costola} \}$

- Show the Front-coding compression of S
- Show the Locality-preserving front-coding compression of S, setting $c=1$ and counting (for simplicity) previous characters and digits for the “backward copy-check” in the LPFC-algorithm

$\langle 0, \beta A C i \rangle \langle 3, 0 \rangle \langle 1, 0 S S \rangle \langle 1, U O N O \rangle \langle 2, U \rangle \langle 0, C O S T I \rangle \langle 4, \alpha F S \rangle$

Let us given the set of strings $S = \{ \text{aaabbb}, \text{aabdb}, \text{abbdb}, \text{abbbe}, \text{abbc}, \text{abc} \}$

- Describe a two-level index that uses in internal memory a Patricia trie on 2 strings
- Show how it is searched the string “aaabcb” in that 2-level index



$\checkmark \langle 2, 200 \& 8 \rangle, \langle 4, 1 \rangle, \langle 1, 555 \rangle$

$\langle 0, 1000 \rangle, \langle 3, 1 \rangle, \langle 2, 1 \rangle$

Given a sequence of integers $S = (1, 6, 15, 18, 19, 20, 21)$, encode them using (by motivating which sequence is compressed):

- Rice code with $k=3$
- Interpolative coding (the first 3 integers only: the ones in positions 2, 4, 6)

$$x \geq 1 \quad R_k(x) \begin{cases} q = \left\lfloor \frac{x-1}{2^k} \right\rfloor \text{ quotient} \rightarrow \text{Encoded as } U(q+1) \\ r = x - 1 - 2^k \cdot q \text{ remainder} \rightarrow \text{Encoded as } L_{2^k}(r) \text{ } k \text{ bits} \end{cases}$$

Just encoding: 1, 5, 9, 3, 1, 1, 1

$$\Rightarrow k=3 \quad R_3(1) = \begin{cases} q = \left\lfloor \frac{1-1}{8} \right\rfloor = 0 & U(1) = 1 \quad R_3(1) = 1000 \\ r = 1 - 1 - 8 \cdot 0 = 0 & 000 \end{cases}$$

$$R_3(5) = \begin{cases} q = \left\lfloor \frac{5-1}{8} \right\rfloor = 0 & U(1) = 1 \quad R_3(5) = 1100 \\ r = 5 - 1 - 8 \cdot 0 = 4 & 000 \end{cases}$$

$$R_3(9) = \begin{cases} q = \left\lfloor \frac{9-1}{8} \right\rfloor = 1 & U(2) = 01 \quad R_3(9) = 01000 \\ r = 9 - 1 - 8 \cdot 1 = 0 & 000 \end{cases}$$

$$R_3(3) = \begin{cases} q = \left\lfloor \frac{3-1}{8} \right\rfloor = 0 & U(1) = 1 \quad R_3(3) = 1010 \\ r = 3 - 1 - 8 \cdot 0 = 2 & 010 \end{cases}$$

1, 6, 15, 18, 19, 20, 21

$l=1, r=7, low=1, high=21$

$$m = \left\lfloor \frac{l+r}{2} \right\rfloor = 4 \quad S_m = 18$$

$$a = low + m - l = 1 + 4 - 1 = 4$$

$$b = high - (r-m) = 21 - (7-4) = 18$$

$S_m - a = 19$ and if we $\lceil \log_2 b - a + 1 \rceil = 6$ but $(1110)_2$

$<1, 3, 1, 17>$

$<5, 7, 19, 21>$

$$m = \left\lfloor \frac{1+3}{2} \right\rfloor = 2 \quad S_m = 6$$

$$m = \left\lfloor \frac{5+7}{2} \right\rfloor = 6 \quad S_m = 20$$

$$a = 1 + 2 - 1 = 2$$

$$a = 19 + 6 - 5 = 20$$

$$b = 17 - (3-2) = 16$$

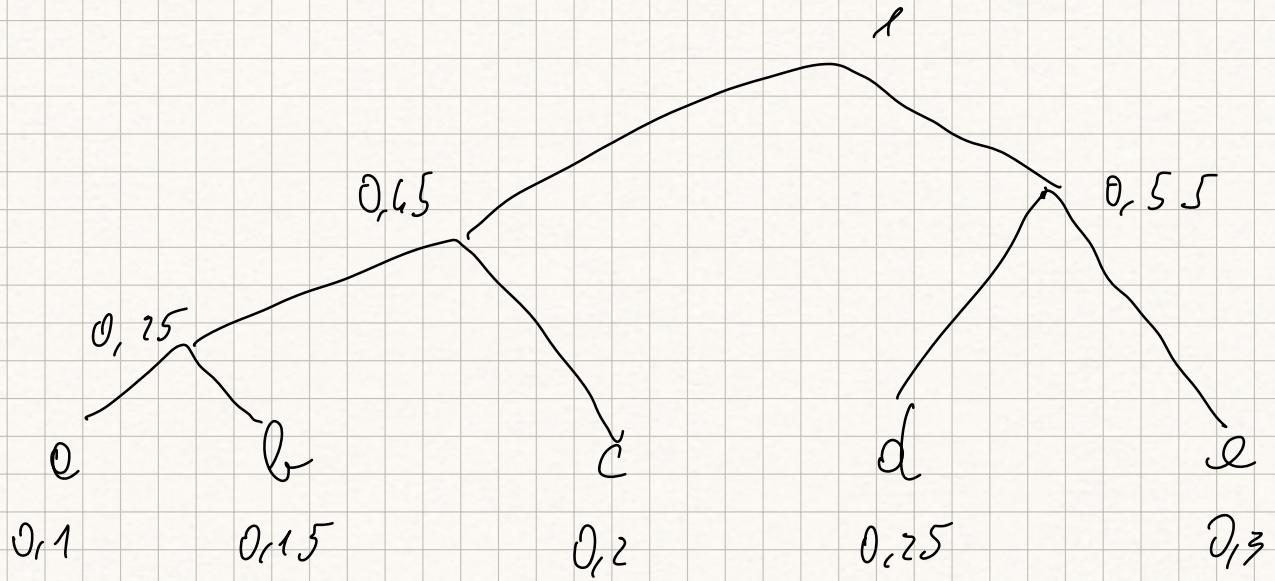
$$b = 21 - (7-6) = 20$$

$$S_m - a = 4 \quad \text{and } \lceil \log_2 (16-2+1) \rceil = 4 \quad (0100)_2$$

our bit entered

Given the probabilities $p(a)=0.1$; $p(b)=0.15$; $p(c)=0.2$; $p(d)=0.25$; $p(e)=0.3$,

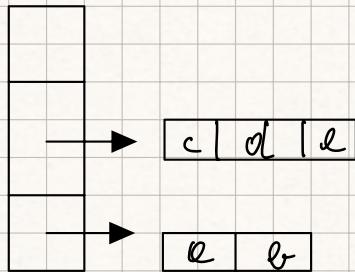
1. Construct the Canonical Huffman code, showing the steps followed by the algorithm.
2. Then use it to decode the bit sequence 1001001, showing each decoding step.



$$\text{num}[1, 3] = [0, 3, 2]$$

$$FC[3] = 0$$

$$\text{Msub}(1, 3) =$$



$$FC[i] = \frac{FC[i+1] + \text{num}[i+1]}{2} =$$

$$FC[2] = \frac{0 + 2}{2} = 1$$

$$FC[1] = \frac{1 + 3}{2} = 2 \quad \text{we have no else in first level}$$

~~1 0 | 0 0 1 | 0 0 1~~

$$v = 1 \quad l = 4 \quad v < FC[e] \\ 1 < 2$$

$$v = 2 \cdot v + \text{nextbit}()$$

$$l \leftarrow l +$$

$$v = 2 + 0$$

$$l = 2$$

~~$v = 0 \quad l = 1 \quad 1 < 2$
 $v = 1 \quad l = 2 \quad 1 < 1$~~

$$2 < 1$$

$$\text{Msub}[l, v - FC[e]] = d$$

$$\text{Msub}[l, v - FC[e]] = c$$

$v=0 \quad l=1 \quad 0 \leq 2$

$v=2 \quad l=2 \quad 0 \leq 1$

$v=1 \quad l=3 \quad 1 \leq 0 \quad \text{symbol}[l, v - \text{PC}[l]] = b$

$\begin{matrix} 3 \\ | \end{matrix} \quad 1$

Given the sequence of integers $S = 1 2 3 4 6 10 11$,

- compress the entire S via Interpolative Coding
- compress the gaps between the consecutive integers of S by PForDelta using $b=2$ bits, and base = 1.

$\langle l=1, r=7, low=1, high=11 \rangle$

$m = \left\lfloor \frac{l+2}{2} \right\rfloor = 4 \quad S_m = 4$

$a = low + (m-l) = 1 + 4 - 1 = 4 \quad a \leq b \leq 8 \quad S_m - a = 4 - 4 = 0$

$b = high - (r-m) = 11 - (7-4) = 8 \quad m \lceil \log_2(b-a+1) \rceil = 3 \quad 000$

$\langle l=1, r=3, low=1, high=3 \rangle \quad \langle l=5, r=7, low=5, high=11 \rangle$

$m = 2 \quad S_m = 2$

$m = \frac{l^2}{2} = 6 \quad S_m = 10$

$a = 1 + 2 - 1 = 2$

$a = 5 + (6 - 5) = 6$

$b = 3 - (3 - 2) = 2$

no division

$b = 11 - (7 - 6) = 10$

$S_m - a = 10 - 6 = 4 \quad m \lceil \log_2 5 \rceil = 3$

100

$\langle l=2=5, low=5, high=9 \rangle$

$\langle l=1=7, low=11, high=11 \rangle$

$$m = 5 \quad S_m = 6$$

$$m = 7 \quad S_m = 11$$

$$Q = 5 + (5 - 5) = 5$$

$$Q = 11 + (7 - 7) = 11$$

new bit emitted

$$B = 9 - (5 - 5) = 9$$

$$L = 11 - (7 - 7) = 11$$

$$S_m - Q = 1 \quad m \quad \text{Flag}_2[5] = 3$$

001

~~11111241
0000130~~

$$B = 2 \quad \text{here} = 1$$

00 00 00 00 01 11 00

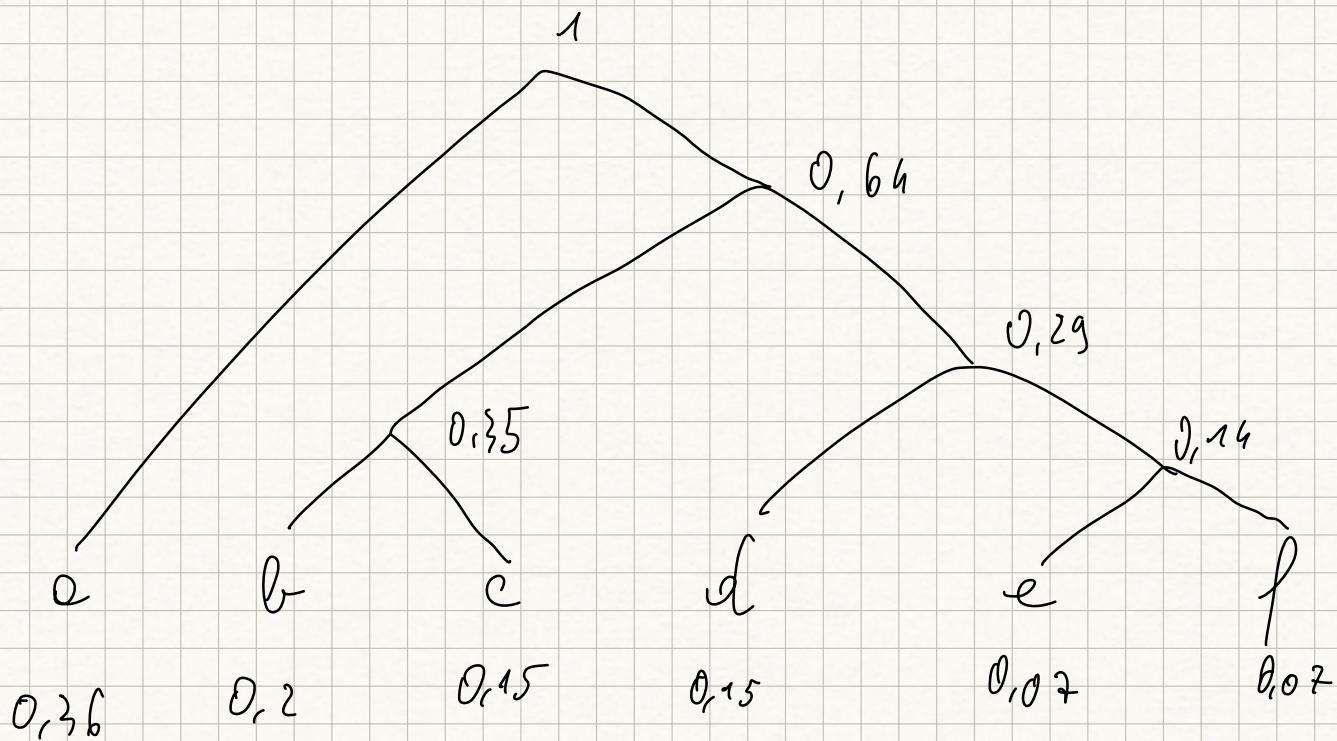
00 $\rightarrow 0$

01 $\rightarrow 1$

10 $\rightarrow 2$

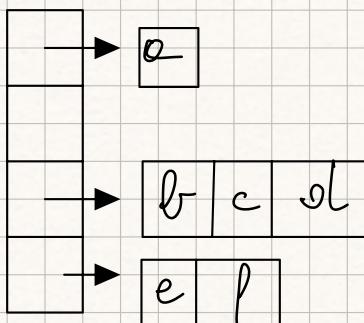
11 $\rightarrow \text{exp}$

Let us given symbols and probabilities: $p(f) = p(e) = 0.07$, $p(c)=p(d) = 0.15$, $p(b)=0.2$, $p(a)=0.36$. Construct the Canonical Huffman code, showing how it is obtained step-by-step.



$$\text{num } [1, h] = [1, 0 \ 3, 2]$$

symbol =



$$FC[0] = 0$$

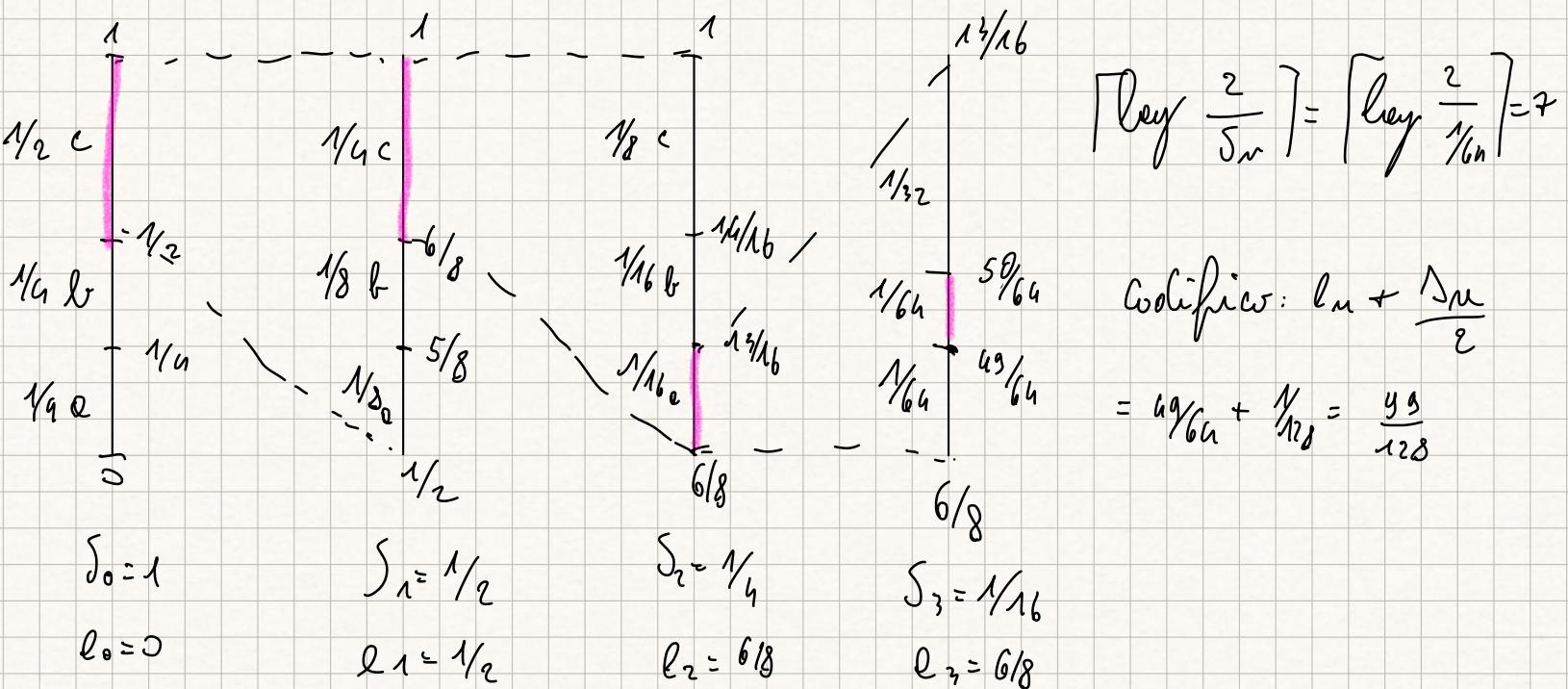
$$FC[3] = \frac{0 + 2}{2} = 1$$

$$FC[2] = \frac{1 + 3}{2} = 2$$

$$FC[1] = \frac{2 + 0}{2} = 1$$

Given the text $T = ccab$, compress it via Arithmetic coding, using as probabilities of the symbols their frequency of occurrence in T .

$$P(0) = 1/4 \quad P(b) = 1/4 \quad P(c) = 1/2$$



$$2 \cdot \frac{99}{128} = \frac{99}{64} \geq 1 \quad \text{exit 1}, \quad \frac{99}{64} - 1 = \frac{35}{64}$$

$$2 \cdot \frac{35}{64} = \frac{70}{64} \geq 1 \quad \text{exit 1}, \quad \frac{70}{64} - 1 = \frac{8}{64} = \frac{1}{8}$$

$$2 \cdot \frac{7}{32} = \frac{6}{32} < 0 \quad \text{ent } 0,$$

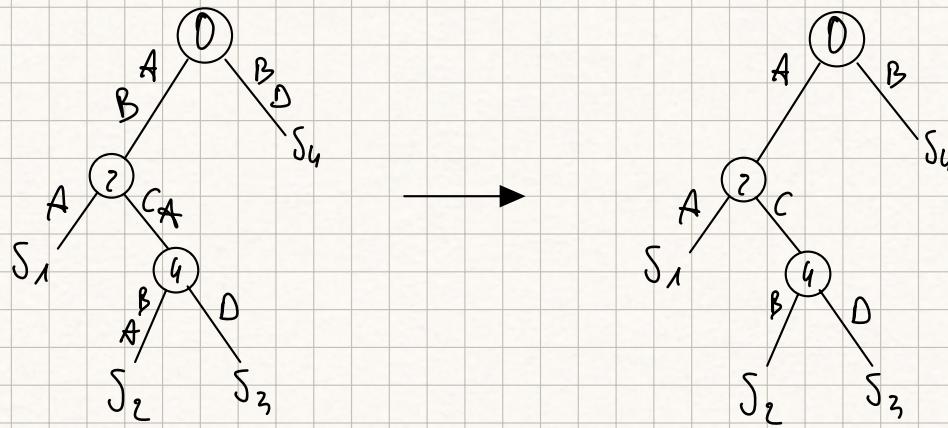
$$2 \cdot \frac{6}{32} = \frac{12}{32} < 0 \quad \text{ent } 0,$$

$$2 \cdot \frac{12}{32} = \frac{24}{32} < 0 \quad \text{ent } 0,$$

$$2 \cdot \frac{24}{32} = \frac{48}{32} \geq 1 \quad \text{ent } 1 \quad \frac{48}{32} - 1 = \frac{16}{32} = \frac{1}{2}$$

$$2 \cdot \frac{1}{2} = 1 \geq 1 \quad \text{ent } 1$$

Build a Patricia trie on the set of strings $S=\{\text{aba}, \text{abcaba}, \text{abcd}, \text{bd}\}$ and show how it is executed a lexicographic search for the string $P_1 = \text{acdac}$, and $P_2 = \text{aab}$



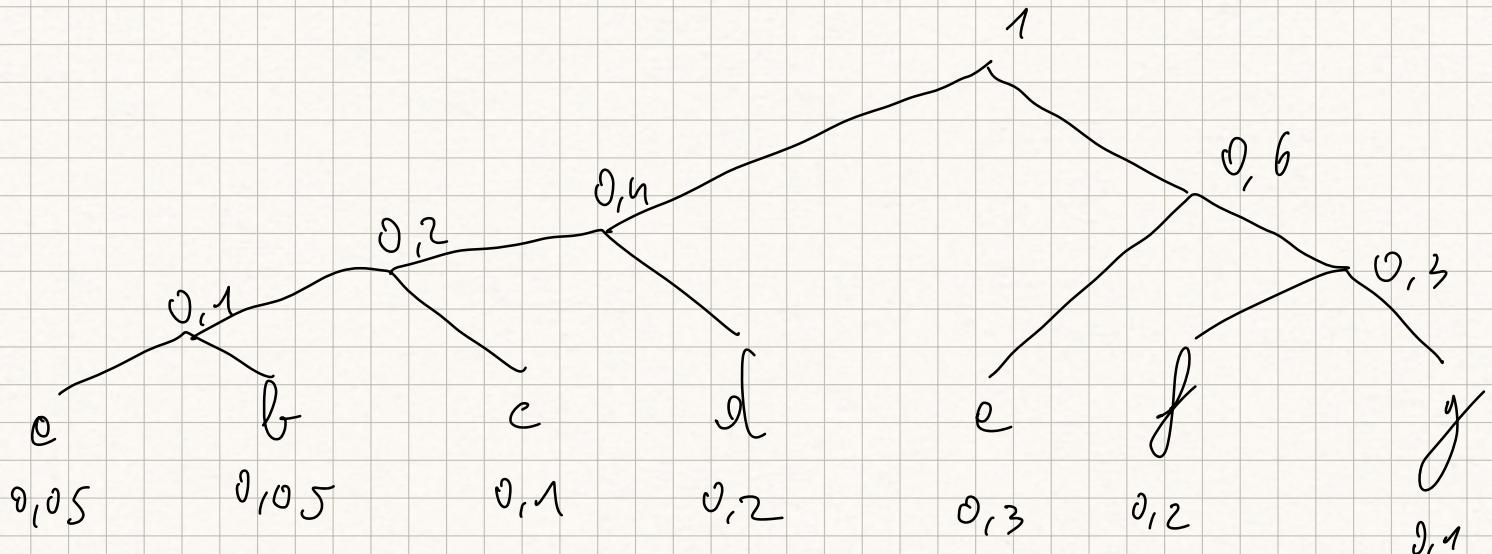
$P_1 = \underset{0}{\text{A}} \underset{1}{\text{C}} \underset{2}{\text{D}} \underset{3}{\text{A}} \underset{4}{\text{C}}$ (cp the (S_1, P_1)) = A cerew the $S_2 \text{ e } S_4$

$P_2 = \underset{0}{\text{A}} \underset{1}{\text{A}} \underset{2}{\text{B}}$ (cp the (S_1, P_2)) = A cerew prime oh S_1

Let us given the symbols and their probabilities: $p(a) = p(b) = 0.05$, $p(c) = p(g) = 0.1$, $p(d) = p(f) = 0.2$, $p(e) = 0.3$.

- Compute the Canonical Huffman code for this distribution
- Decode the first 3 symbols of the coded sequence

10001011011100...



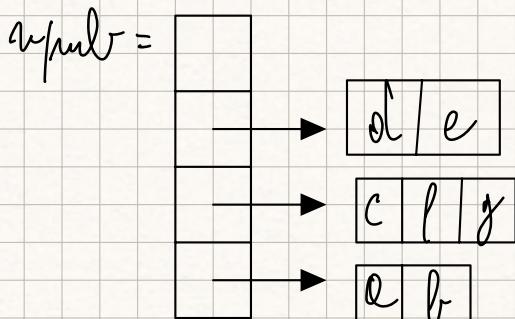
$$H_{\text{min}}(1, 6) = [0, 1, 2, 3, 2]$$

$$FC[4] = 0$$

$$FC[3] = \frac{0+2}{2} = 1$$

$$FC[2] = \frac{1+3}{2} = 2$$

$$FC[1] = \frac{2+2}{2} = 2$$



1 0 | 0 0 0 | 0 1 1 0 1 1 1 0 0

$$v=1 \quad l=1 \quad 1 < 2$$

$$v=2+0 \quad l=2 \quad 2 \geq 2$$

$$\text{symbol}[2, 2-FC[2]] = \textcircled{d}$$

$$v=0 \quad l=1 \quad 0 < 2$$

$$v=0 \quad l=2 \quad 0 < 2$$

$$l=1 \quad l=1 \quad 1 < 1$$

$$v=1 \quad l=3 \quad i \geq 1$$

$$m_{\text{mult}}[3, 1-1] = \textcircled{c}$$

$$v=0 \quad l=1 \quad 0<2$$

$$v=1 \quad l=2 \quad 1 < 2$$

$$v=2 \quad l=3 \quad 2 > 1$$

$$m_{\text{mult}}[3, 3-1] = \textcircled{y}$$

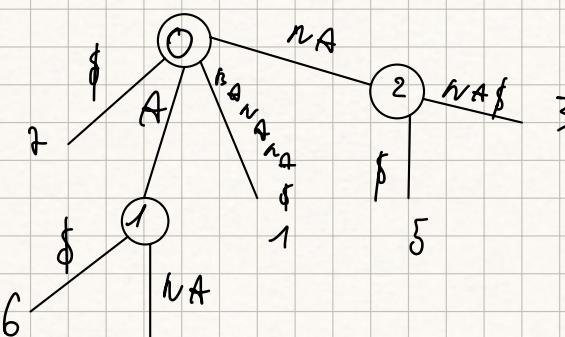
Given the string $S=baobab$, compute the Suffix Array, and then derive the LCP array by using the linear-time algorithm seen in class, showing the amount of characters really compared when determining the various LCPs.

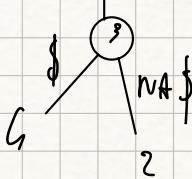
$BAOBAB\$ \rightarrow \begin{matrix} BAODA B \$ \\ 1 \end{matrix}, \begin{matrix} AOBAB \$ \\ 2 \end{matrix}, \begin{matrix} OBA B \$ \\ 3 \end{matrix}, \begin{matrix} BAB \$ \\ 4 \end{matrix}, \begin{matrix} AB \$ \\ 5 \end{matrix}, \begin{matrix} B \$ \\ 6 \end{matrix}, \begin{matrix} \$ \\ 7 \end{matrix}$

ST	SUFF
0	\$
1	AB\$
2	AOBAB\$
3	B\$
4	DA B \$
5	BAOBA B \$
6	OBAB\$
7	

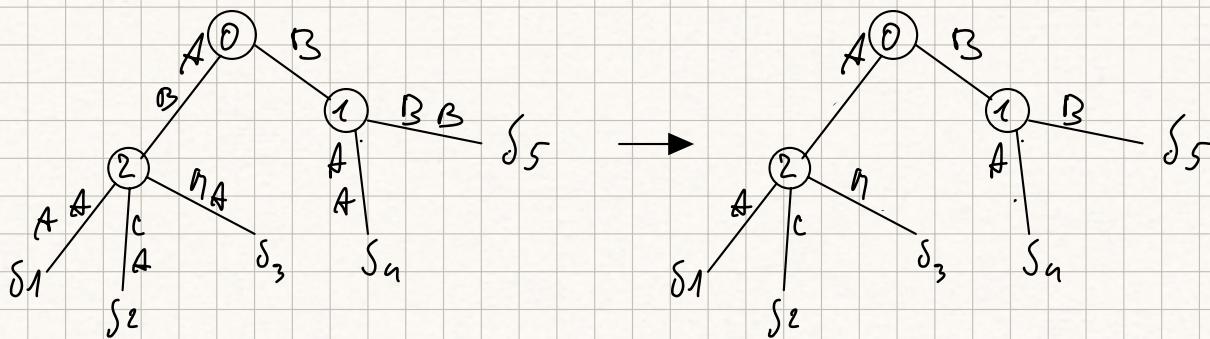
1. Build the suffix tree of the string $T = BANANA\$$.

$BANANA\$$, $ANANA\$$, $NANA\$$, $NA\$$, $A\$$, $\$\$$, $\$$





Given the set of strings $S = \{abaa, abca, abma, baa, bbb\}$, build a Patricia trie and show the steps for the lexicographic search of the strings $P_1 = aaa$, $P_2 = abb$



$P_1 = \begin{matrix} A & A & A \\ \uparrow & \uparrow & \uparrow \\ 0 & 1 & 2 \end{matrix}$ prenow S_1 $LCP(S_1, P_1) = A$ quiandi cercw prime iki S_1

$P_2 = \begin{matrix} A & B & B \\ \uparrow & \uparrow & \uparrow \\ 0 & 1 & 2 \end{matrix}$ ho mismatch in position 3, prenow quinshi strings notif il nowb ②

$LCP(S_1, P_2) = 2 = AB$ quiandi cercw tre $S_1 S_2$

Given the sequence of integers $(11, 14, 16, 17, 19, 20, 21, 31)$, show how to encode them based on

- Elias-Fano Code
- PForDelta with base = 11, and $b=3$

$$m=32 \quad b = \lceil \log_2 m \rceil = 5 \quad l = \lceil \log_2 \frac{m}{n} \rceil = 2 \quad h = b - l = 5 - 2 = 3$$

01011	11	$L = 111000011100011$
01110	14	
10000	16	$M = 00101011101100110$
10001	17	10

100	11	19
10	100	20
10	101	21
11	111	22
h	l	

bore = 1 b=3

∴ : 0, 3, 5, 6, 8, 9, 10, 20

gap : 0, 3, 2, 1, 2, 1, 1, 10

000 011 010 001 010 001 001 111

000	0
001	1
010	2
011	3
100	4
101	5
110	6
111	7 exp.

Given the sequence of integers (11, 14, 16, 19, 20, 21, 22), show how to encode them based on

- Elias-Fano Code
- Interpolative Code (just one level of recursion, hence 3 numbers)

$$\mu = 23 \quad b = \lceil \log_2 \mu \rceil = 5 \quad l = \lceil \log_2 \frac{\mu}{m} \rceil = 4 \quad h = 5 - 4 = 1$$

01011	14
01110	19
10000	16
10011	19
10100	20
10101	21
10110	22

L = 1101001100011010

H = 001010110111000

$\langle l=1, r=8, low=19, high=22 \rangle$

$$m = \left\lceil \frac{l+r}{2} \right\rceil = 4 \quad S_m = 19$$

$$Q = low + (m - l) = 19 + (4 - 1) = 22$$

$$S_m - Q = 5 \quad m \lceil \log_2 (22 - 19 + 1) \rceil = 3 \\ (101)_2$$

$$b = high - (r - m) = 22 - (8 - 4) = 16$$

$\langle l=1, r=3, low=11, high=18 \rangle$

$$m = \left\lfloor \frac{l+r}{2} \right\rfloor = 2 \quad S_m = 14$$

$$Q = 14 + (2 - 1) = 15$$

$$S_m - Q = 2 \quad m \lceil \log_2 (18 - 14 + 1) \rceil = 3$$

$$b = 18 - (3 - 2) = 17$$

$$(010)_2$$

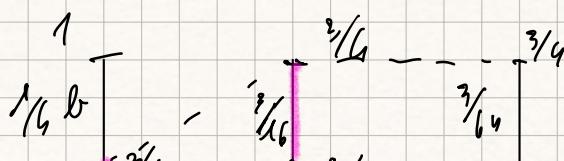
$\langle l=5, r=7, low=20, high=22 \rangle$

$$m = \frac{l+r}{2} = 6 \quad S_m = 21$$

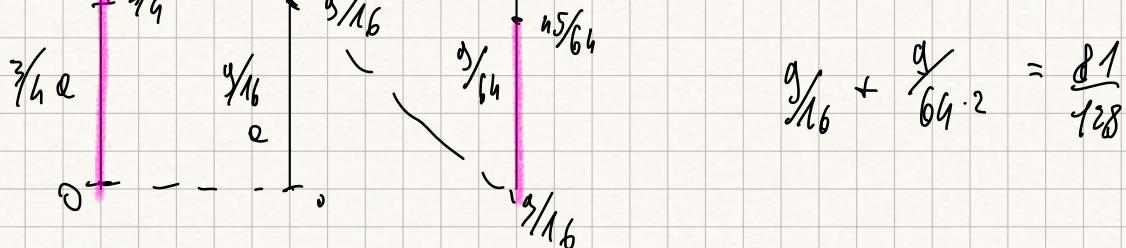
$$Q = 20 + (6 - 5) = 21 \quad \text{no bit omitted}$$

$$b = 22 - (7 - 6) = 21$$

Given the string T = aba, and the probabilities $P(a) = \frac{3}{4}$ and $P(b) = \frac{1}{4}$, show the result of the Arithmetic Coding applied on T. (hint: Please work with the dyadic fractions, not change them into reals.)



$$d = \lceil \log_2 \frac{2}{\frac{3}{16}} \rceil = 4$$



$$S_0 = 1 \quad S_1 = 3/4 \quad S_2 = 5/6$$

$$l_0 = 0 \quad l_1 = 0 \quad l_2 = 9/16$$

$$\frac{9}{16} + \frac{9}{64} \cdot 2 = \frac{81}{128}$$

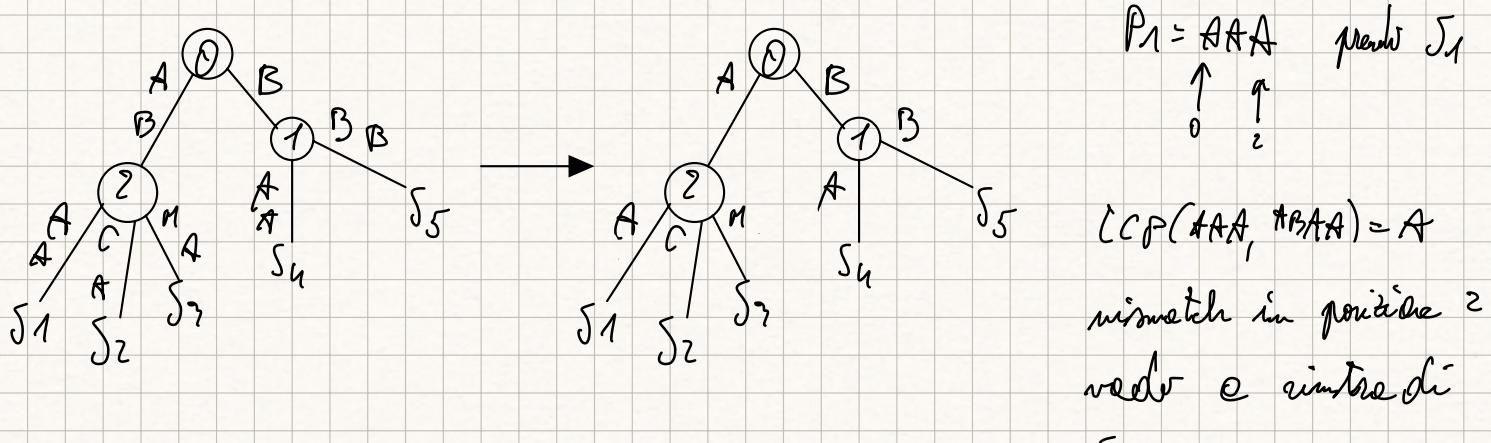
$$\frac{81}{128} \cdot 2 = \frac{81}{64} > 1 \text{ exit } 1 \quad \frac{81}{64} \cdot 1 = \frac{17}{64}$$

$$\frac{17}{64} \cdot 2 = \frac{17}{32} < 1 \text{ exit } 0$$

$$\frac{17}{32} \cdot 2 = \frac{17}{16} > 1 \text{ exit } 1 \quad \frac{17}{16} - 1 = \frac{1}{16}$$

$$\frac{1}{16} \cdot 2 = \frac{1}{8} < 1 \text{ exit } 0$$

Given the set of strings $S = \{\text{abaa}, \text{abca}, \text{abma}, \text{baa}, \text{bbb}\}$, build a Patricia trie and show the steps for the lexicographic search of the strings $P_1 = \text{aaa}$, $P_2 = \text{abb}$.



$P_1 = \text{AAA}$ prefix S_1

$P_2 = \text{ABB}$

↑ ↑
0 1
mismatch, prefix position before all others ② $\text{LCP}(S_1, P_2) = AB$

cover the S_1 e S_2

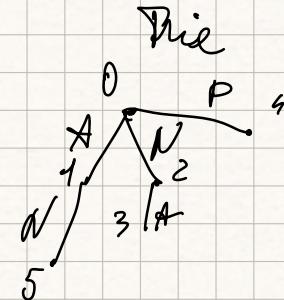
Given the string $S = \text{annapanna}$, compute its parsing LZ77 and

ANNA PANNA

 $\langle 0, 0, A \rangle, \langle 0, 0, N \rangle \langle 1, 1, A \rangle \langle 0, 0, P \rangle \langle 5, 4, EOF \rangle$ LZ78

$\langle 1, N \rangle \langle 3, EOF \rangle$

A: N: N: A: P: A: N: N: A
 $\langle 0, A \rangle \langle 0, N \rangle \langle 2, A \rangle \langle 0, P \rangle$



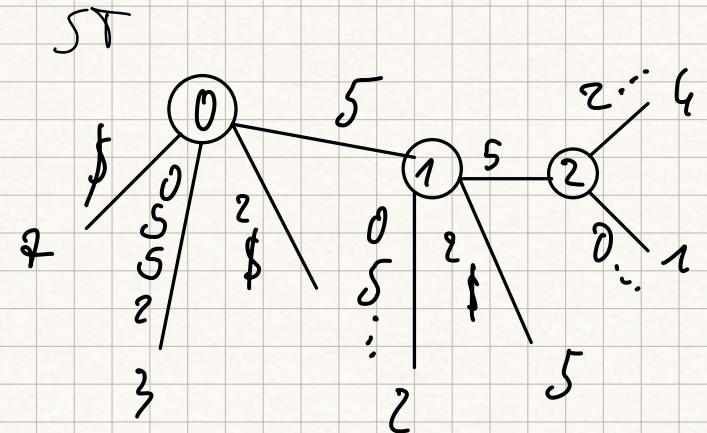
Question #4 [scores 3+3+2]. Given the string T formed by your Matricola, and hence consisting of 6 digits:

- Show the suffix array of T, in which every digit is interpreted as a symbol;

Form the string P as given by the two middle digits of T (i.e. if the Matricola is 123456, then P=34). Then describe the algorithm that efficiently counts the occurrences of P in T.

- Comment on the time complexity of the counting algorithm as a function of n (= T's length), p (= P's length) and the number occ of P's occurrences.

	LCP	ST
1	550 552 \$	0
2	50 552 \$	0
3	0 552 \$	0
4	552 \$	1
5	52 \$	1
6	2 \$	2
7	\$	2



Question #2 [scores 5]. Take your Matricola (of 6 digits), change every occurrence of 0 with 1 (if any), and then interpret each digit as an integer gap, and finally derive an increasing integer sequence by summing those gaps: namely, if the Matricola is 120304, then you transform it into 121314, and then you get the corresponding integer sequence as 1, 3 (=1+2), 7 (=1+2+1), 7 (=1+2+1+3), 8 (=1+2+1+3+1), 12(=1+2+1+3+1+4). Compress the resulting increasing integer sequence with Elias-Fano.

$$551552 \Rightarrow 5, 10, 11, 16, 21, 23$$

$$\mu = 24 \quad \ell = \lceil \log_2 24 \rceil = 5$$

00101	5
01010	10
01011	11
10000	16
10101	21
10111	23

$$\ell = \lceil \log_2 \frac{\mu}{m} \rceil = 2 \quad b = 3$$

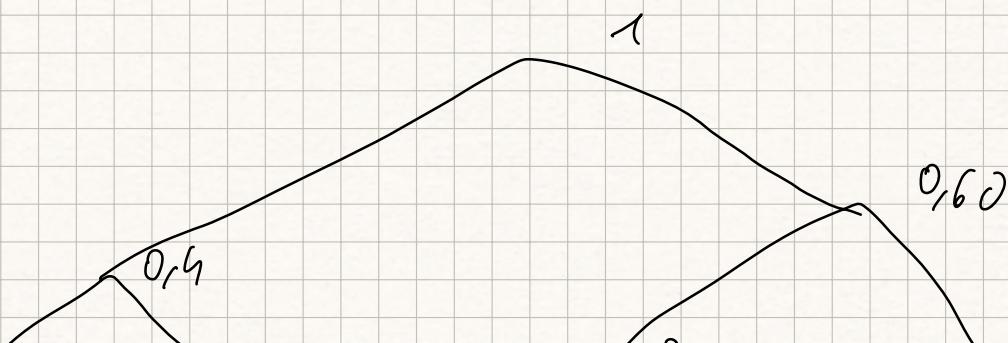
$$\begin{aligned} l &= 011011000111 \\ u &= 010110011011000 \end{aligned}$$

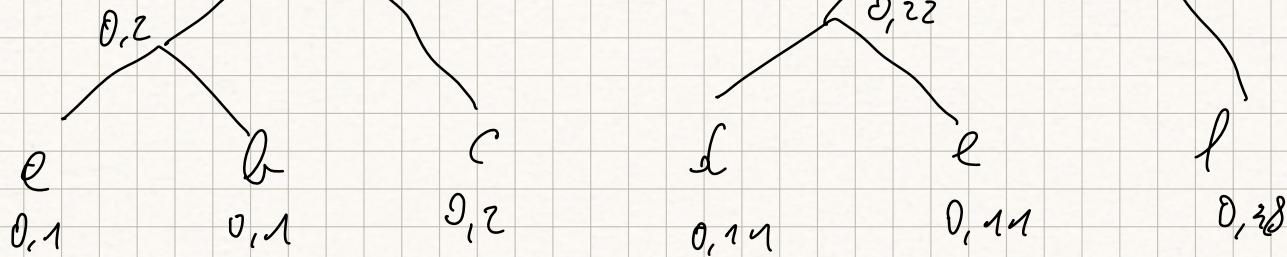
Question #1 [scores 2+3+4]. Given the symbols and their probabilities: $p(a) : p(b) = 0.1$, $p(c) = 0.2$, $p(d)=p(e)=0.11$, $p(f)=0.38$.

Compute the Huffman code for this distribution.

Compute the Canonical variant of the Huffman code (by sorting alphabetically the letters in every SYMB's list).

Decode the first 2 symbols of the coded sequence 11010.

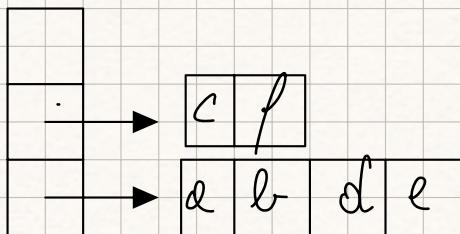




$$\text{num}(1,3) = [0, 2, 6]$$

$$Fc[3] = 0$$

$$\text{ymb}(1,3) =$$



$$Fc[2] = \frac{0+4}{2} = 2$$

$$Fc[1] = \frac{2+2}{2} = 2$$

$$n=1 \quad l=1 \quad 1 < 2$$

$$\text{ymb}[2, 3-2] = f$$

$$n=2 \cdot 1 + 1 \quad l=2 \quad 3 > 2$$

$$n=0 \quad l=1 \quad 0 < 2$$

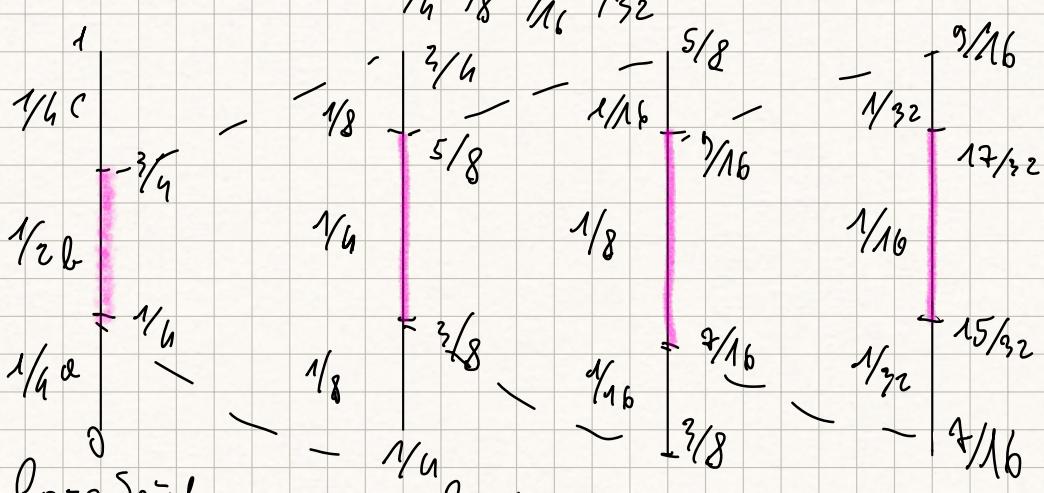
$$n=2 \cdot 0 + 1 \quad l=2 \quad 1 < 2$$

$$\text{ymb}[3, 2] = d$$

$$n=2 \cdot 1 + 0 \quad l=3 \quad 2 > 0$$

Decode the compressed sequence $<4, 011110>$ produced by arithmetic code, by assuming probabilities $P[a]=P[c]=1/4$ and $P[b]=1/2$.

$$0 \quad 1 \quad 1 \quad 1 \quad 1 \quad 0 \quad = \frac{1+2+4+8}{32} = \frac{15}{32}$$



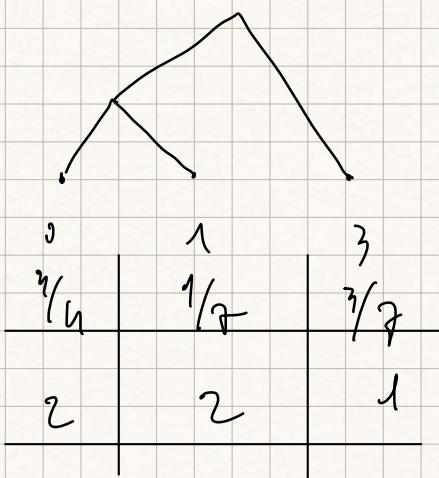
$$Q_1 = \frac{1}{4} S_1 = \frac{1}{2} \quad Q_2 = \frac{3}{8} S_2 = \frac{1}{4} \quad Q_3 = \frac{2}{16} S_3 = \frac{1}{16}$$

Bertrand-Weller

text = AATATA#
 1 2 3 4 5 6 7 8 9 10
 L

5A											① BWT
10	# A										② MTF
9	A #										③ RLE Ø
1	A n										④ Canonical Huffman
7	ATA#										Predictable
5	A T A T A #) position # = 3
3	AT) L = {A, T, }
2	M										1 3 1 1 3 3 1 1 1
8	T A #										A T T T M A A A A → 0 3 1 3 3 0 0
6	T A T A #										
4	T A										

$$P(0) = \frac{3}{7} \quad P(1) = \frac{1}{7} \quad P(2) = \frac{3}{7}$$

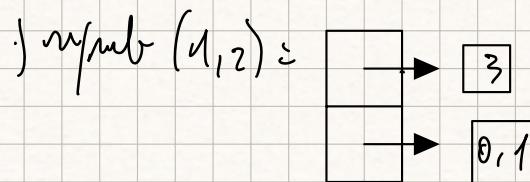


$$L = \{ \overset{1}{A}, \overset{2}{M}, \overset{3}{T} \}$$

$$L = \{ \overset{1}{T}, \overset{2}{A}, \overset{3}{M} \}$$

$$L = \{ \overset{1}{M}, \overset{2}{T}, \overset{3}{A} \}$$

$$\text{num}(1, 2) = 1, 2$$



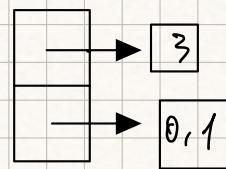
$$\therefore FC[2] = 0$$

$$\therefore FC[1] = \frac{0+2}{2} = 1$$

Body: 00101110000
 0 3 1 3 3 0 0

Decomposition of B-W

$\text{symbol}(1,2) =$



$$B = 00 \mid 101 \mid 110000$$

0 1 3

$$v=0 \quad l=1 \quad 0 < 1$$

$$FC[2] = 0$$

$$v=0 \quad l=2 \quad 0 \geq 0$$

$$FC[1] = 1$$

$$\text{symbol}[2, 0 \leftarrow 0] = 0$$

$$v=1 \quad l=1 \quad 1 \geq 1$$

$$\begin{array}{ccccccccc} 0 & 3 & 1 & 3 & 3 & 3 & 0 & 0 \\ \hline 10 & 11 & & & & & 100 \\ 2 & & 3 & -1 & & & 4 & -1 \\ 1 & -1 & 2 & & & & 3 & \\ \hline \end{array}$$

$$\text{symbol}[1, 0] = 3$$

$$v=0 \quad l=1 \quad 0 < 1$$

$$\begin{array}{cccccccc} 1 & 3 & 11 & 3 & 3 & 1 & 1 & 1 \\ \hline A & T & T & T & M & A & A & A \end{array}$$

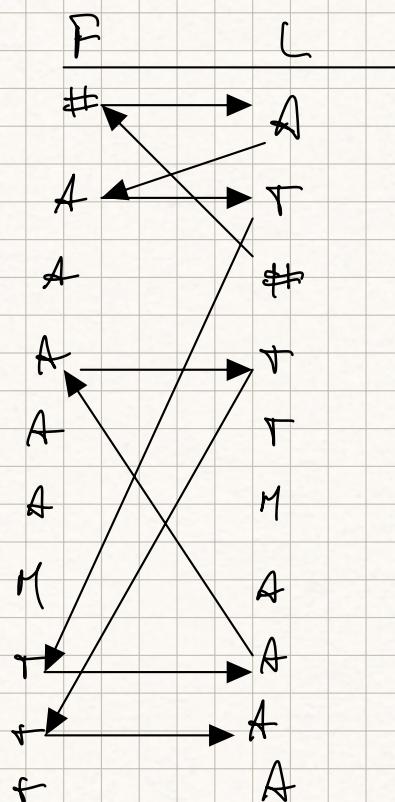
$$\text{symbol}[2, 1] = 1$$

$$\mathcal{L} = (A, M, T)$$

$$\mathcal{L} = (T, A, M)$$

$$\mathcal{L} = (M, T, A)$$

$$\mathcal{L} = (A, M, T)$$



$$T = A \cap A \cap A \cap A \cap A \#$$

A = 000 000 0

B = 000 001 0

C = 000 110 0

D = 000 111 0

E = 100

F = 101 0

2 level-index, in which $B=2$ string

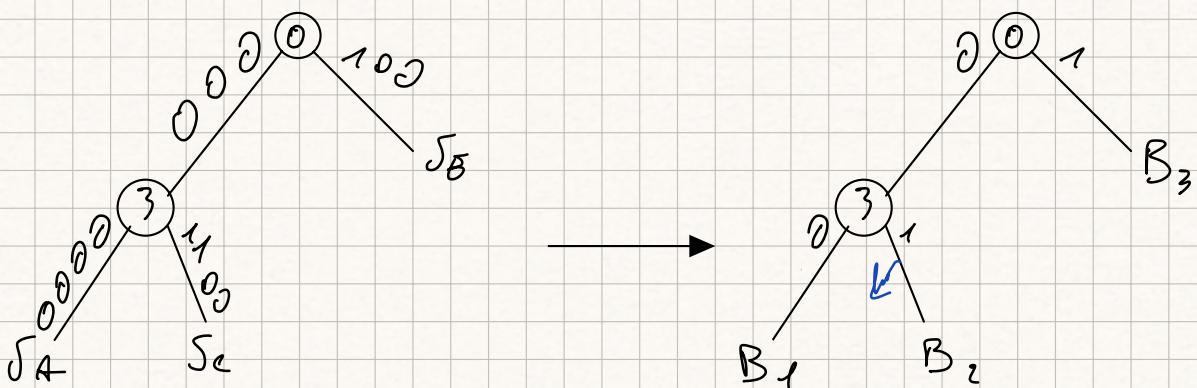
Patricia Trie in int. memory level

FC on disk

$$B_1 = \boxed{<0, 0000000> <5, 10>}$$

$$B_2 = \boxed{<0, 000 110 0> <5, 10>}$$

$$B_3 = \boxed{<0, 100> <2, 10>}$$



Lexicographic search $P = 000 \ 1\boxed{0}1$ $lcp(P, C) = \frac{0001}{4}$

\uparrow
 0
 \uparrow
 3
↑ mismatch

We search in B_1 .

$P_1 = 10\#$ $lcp(P_1, B) = 2$ jumps in B_3 and scan

$P_2 = 10\#\#$ $lcp(P_2, E) = 2$ jumps in B_3 and scan

If you are to the left you have to scan the previous block, if you are to the right you have to scan the right block

Given the string $S = abababc$ show the result of the algorithmic pipeline BWT + MTF + RLE0 + Huffman, where RLE0 is the application of the special RunLengthEncoding algorithm over 0-runs and the Wheeler code.

$A \ B \ A \ B \ A \ B \ C \ \#$

$\# \ A \ B \ A \ B \ A \ B \ C$

$C \ \# \ A \ B \ A \ B \ A \ B$

$B \ C \ \# \ A \ B \ A \ B \ A$

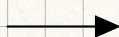
$A \ B \ C \ \# \ A \ B \ A \ B$

$B \ A \ B \ C \ \# \ A \ B \ A$

$A \ B \ A \ B \ C \ \# \ A \ B$

$B \ A \ B \ A \ B \ C \ \# \ A$

Sorting



F

$\# \ A$

$A \ B$

$A \ B$

$A \ B$

$B \ A$

$B \ C$

$C \ \#$

L

C

$\#$

B

B

A

A

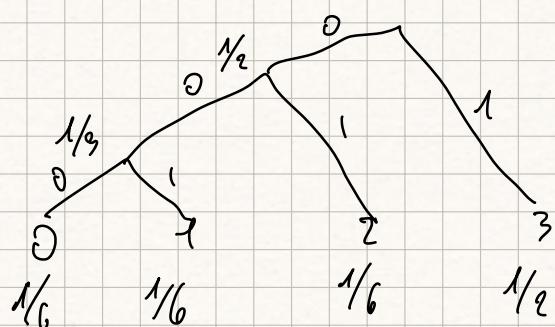
B

Preamble

in pos. 2

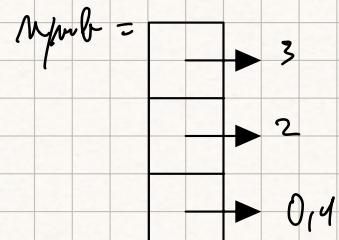
$$P(0) = \frac{1}{6} \quad P(1) = \frac{1}{6}$$

$$P(2) = \frac{1}{6} \quad P(3) = \frac{1}{2}$$



$$\text{Body} = 3 \ 3 \ 0 \ 3 \ 1 \ 2$$

$$\text{num}(1, 3) = (1, 1, 2)$$



$$FC[3] = 0$$

$$FC[2] = \frac{0+2}{2} = 1$$

$$FC[1] = \frac{1+1}{2} = 1$$

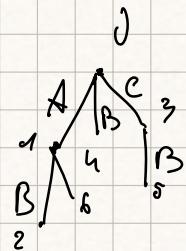
Question #1 [scores 4+6+6] Given the string $S = aabcbcbada$, show the:

- LZ77 parsing of S.
- LZ78 parsing of S, along with the auxiliary data structure used to compute it.
- Canonical Huffman encoding of S

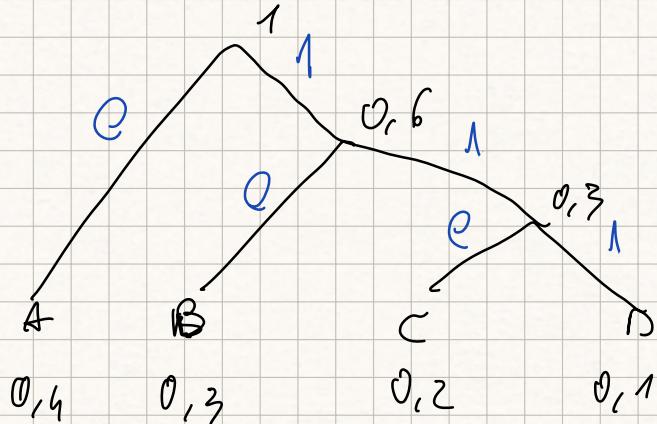
$\begin{array}{c} 1 \\ A | A \\ \hline 2 \\ B | C | \end{array} \begin{array}{c} 3 \\ B | C \\ \hline 4 \\ B | A | D | A \end{array}$

1) $\langle 0, 0, A \rangle \langle 1, 1, B \rangle \langle 0, 0, C \rangle \langle 2, 3, 4 \rangle \langle 0, 0, D \rangle \langle 2, 1, 0, P \rangle$

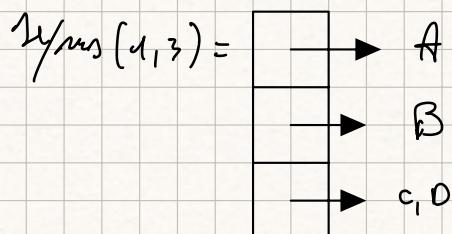
2) $\langle 0, A \rangle, \langle 1, B \rangle, \langle 0, C \rangle, \langle 0, B \rangle, \langle 3, B \rangle, \langle 1, P \rangle, \langle 1, 0, P \rangle$



$$P(A) = 6/10 \quad P(B) = 3/10 \quad P(C) = 3/10 \quad P(D) = 1/10$$



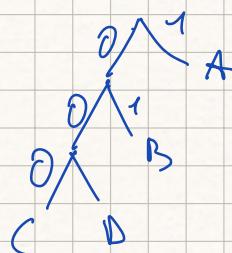
$$\text{mme}(1, 3) = (1, 1, 2)$$



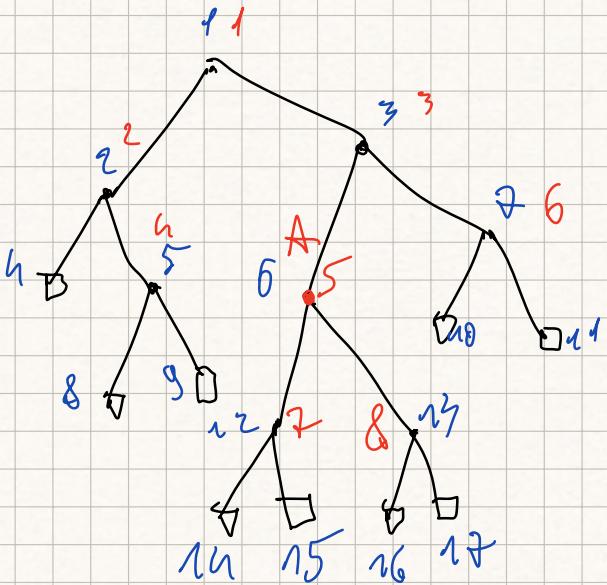
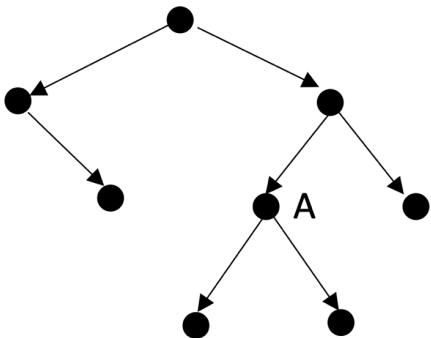
$$FC[3] = 0$$

$$FC[2] = \frac{0+2}{2} = 1$$

$$FC[1] = \frac{1+1}{2} = 1$$



Question #4 [scores 4] Show the succinct representation of the following binary tree. Then show how to use this representation to navigate from the root to the node labelled A, and then back to the parent of A.



$$B = \begin{matrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \end{matrix}$$

Rank_{k+1}[2x] → go left if $B[2x] = 1$

Rank_{k+1}[2x+1] → go right if $B[2x+1] = 1$

Rank₁[2·1+1] = 3 $B[3] = 1$ ok

Rank₁[2·3] = 5 $B[5] = 1$ ok

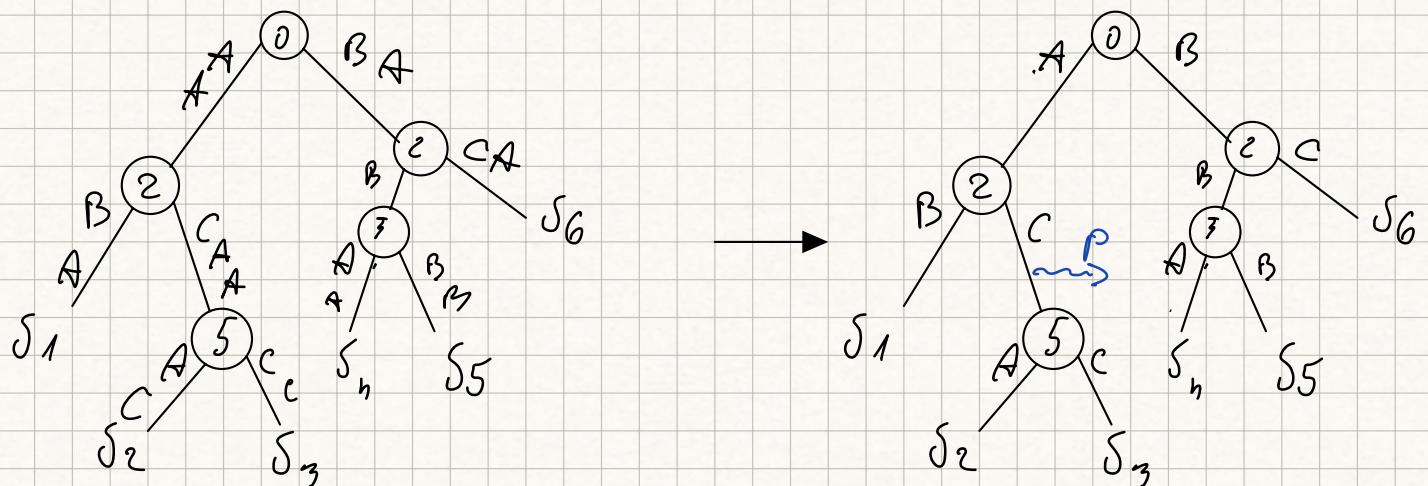
going to parent $\left\lfloor \frac{\text{Select}_1(5)}{2} \right\rfloor = 3$

Given the ordered set of strings:

$$S = \{ AABA, AACAAAC, AACAAACC, BABAA, BABBB, BACA \}$$

- Build the Patricia trie PT for S.

- Show the steps executed to lexicographically search for the pattern P = AACBACD in the set S by means of PT.



$$P = \begin{matrix} A & A & C & B & A & C & D \\ \uparrow & \uparrow & & \uparrow & & & \\ 0 & 2 & & 5 & & & \end{matrix} \quad \text{prefix } S_3 \quad \text{lcp}(P, S_3) = \frac{AAC}{3} \quad \text{cover tree } S_3, S_4$$

Given the string "ABABABC" compress it by using the pipeline BWT + MTF + RLE0 + Huffman, where MTF counts letter's positions from 0, and RLE0 uses the Wheeler's code.

$$T = ABABAABC\$$$

$$\$ A B A B A B C$$

$$C \$ A B A B A B$$

$$B C \$ A B A B A$$

$$A B C \$ A B A B$$

$$B A B C \$ A B A$$

$$A B A B C \$ A$$

$$B A B A B C \$ A$$

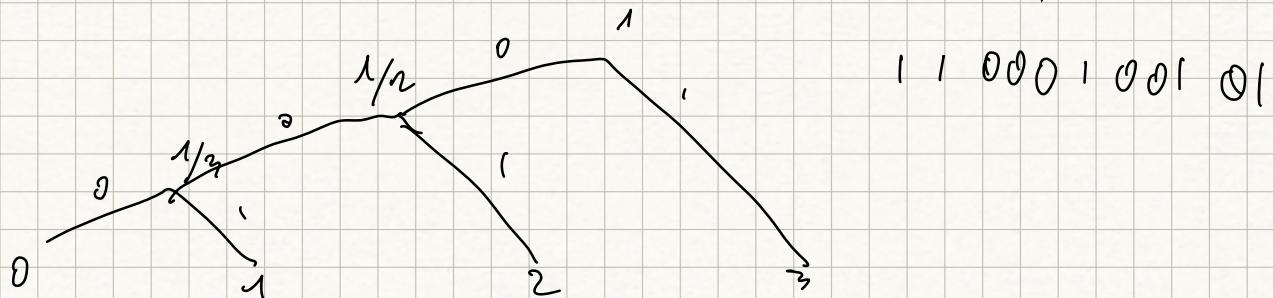
P	L
\\$ A	C
A B	\\$ # pos 2
A B	B
A B	B
B A	A
B A	A
B C	A
C \\$	B

$\begin{matrix} 3 & 3 & 1 & 3 & 1 & 1 & 2 \\ C & B & A & A & A & B \end{matrix}$
 ↓↓↑↑↑↑

$$\mathcal{L} = \{A(B, C)\} \rightarrow \mathcal{L} = \{C, A, B\} \rightarrow \mathcal{L} = \{B, C, A\}$$

3 3 0 3 1 2

$$P(0) = 1/6 \quad P(1) = 1/6 \quad P(2) = 1/6 \quad r(\cdot) = 1/2$$



Given the following SYMB and FC arrays for a Canonical Huffman code: SYMB[1] = [], SYMB[2] = [A,B], SYMB[3] = [C,D,E], SYMB[4] = [F,G]

$$FC = [2, 2, 1, 0]$$

Decompress the first 2 letters of the following compressed bit sequence: 000|111...

$$v=0 \quad l=1 \quad 0 < 2$$

$$v=0 \quad l=2 \quad 0 < 2 \quad \text{symb}[0, 1-0] = 6$$

$$v=0 \quad l=3 \quad 0 < 1$$

$$v=1 \quad l=4 \quad 1 \geq 0$$

$$v=1 \quad l=1 \quad 1 \leq 2 \quad \text{symb}[2, 3-2] = B$$

$$v=2 \cdot 1 + 1 = 3 \quad l=2 \quad 3 \geq 2$$

Given the string "CABABCA" compress it by using the pipeline BWT + MTF + RLE0 + Huffman, where MTF counts letter's positions from 0, and RLE0 uses the Wheeler's code.

CABABCAC#

#CABAABCAC

A#CABAABCAC

CABA#CABAABC

B,CABA#CABAABC

A#B,CABA#CABA

B,A,B,CABA#CABA

A#B,CABA#CABA

P
#C

A#

A B

A B

B A

B C

C A

C A

L
A
C

C

B

A

A

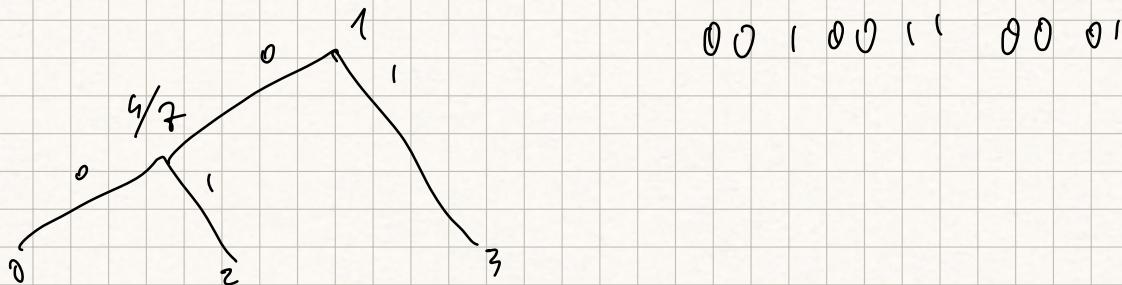
A

B

in pos 8

$$\langle \begin{smallmatrix} 1 & 2 & 1 & 3 & 1 & 2 \\ 8, A & C & C & B & A & A & B \\ 0 & 3 & 0 & 3 & 3 & 0 & 2 \end{smallmatrix} \rangle \quad \mathcal{L} = \{A, B, C\} \leftrightarrow \mathcal{L} = \{C, A, B\} \rightarrow \mathcal{L} = \{B, C, A\}$$

$$P(0) = 3/2 \quad P(1) = 1/2 \quad P(2) = 3/2$$



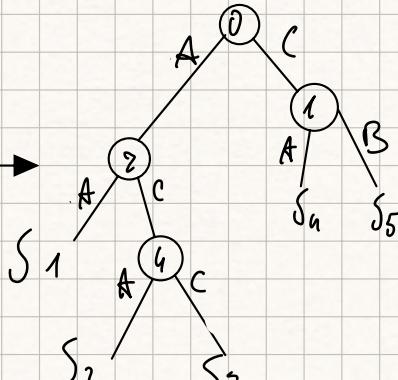
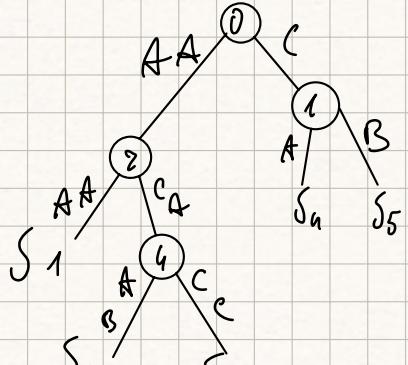
Given the set of strings

$$S = \{\text{AAAA}, \text{AACAAAB}, \text{AACACC}, \text{CA}, \text{CB}\}$$

- Built the Patricia trie PT for S

- Show the steps executed to lexicographically search for P1 =

ABC in the PT



$$P_1 = A B C$$

New S₂

$$\text{lcp}(S_2, P_1) = 1$$

Cerw

tree S₃, S₄

Given the list of integers $A = \{2, 3, 5, 6, 8, 10, 11\}$, use Interpolative Code and emit the encoding of the “first” three integers 6, 3 and 10.

$$l = 1, r = 7, low = 2, high = 11$$

$$m = \left\lfloor \frac{l+2}{2} \right\rfloor = 4 \quad S_m = 6$$

$$Q = 2 + (4 - 1) = 5$$

$$B = 11 - (7 - 4) = 8$$

$$S_{m-Q} = 1 \quad \text{in } \lceil \log_2(8 - 5 + 1) \rceil = 2 \\ (01)_2$$

$$l = 1, r = 3, low = 2, high = 5$$

$$m = 2 \quad S_m = 3$$

$$Q = 2 + (2 - 1) = 3$$

$$3 - 3 = 0 \quad m \lceil \log_2(5 - 3 + 1) \rceil = 1 \\ (0)$$

$$B = 5 - (3 - 2) = 6$$

$$l = 5, r = 7, low = 7, high = 11$$

$$m = 6 \quad S_m = 10$$

$$Q = 7 + (6 - 5) = 8$$

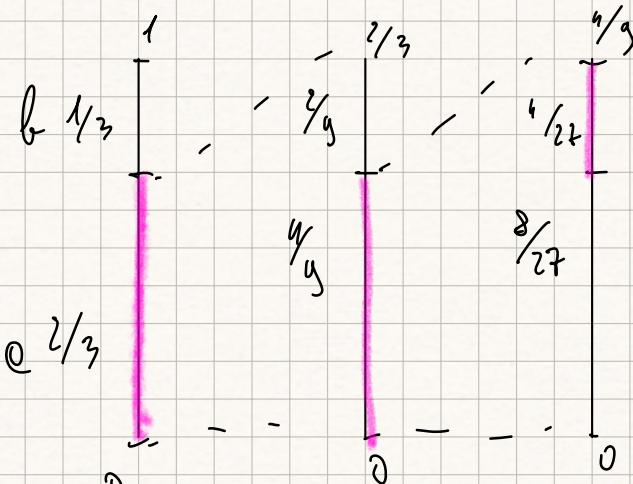
$$S_{m-Q} = 2 \quad m \lceil \log_2(10 - 8 + 1) \rceil = 2 \\ (10)_2$$

$$B = 11 - (7 - 6) = 10$$

Given the text T=aab, encode it using Arithmetic Coding and the

semi-static model. (Associate the intervals from the bottom to the top as "a", "b"; and use the fractions directly.)

$$P(a) = \frac{2}{3} \quad P(b) = \frac{1}{3}$$



$$l_0 = 0 \quad S_0 = 1$$

$$l_1 = 0 \quad S_1 = \frac{2}{3}$$

$$S_2 = \frac{4}{9}$$

$$l_2 = 0$$

$$\frac{S_1}{2} + l_2 = \frac{10}{27}$$

$$\left[\overline{ly}_2, \frac{2}{S_1} \right] = \left[\overline{ly}_2, \frac{27}{2} \right] = 1$$

$$S_3 = \frac{4}{27}$$

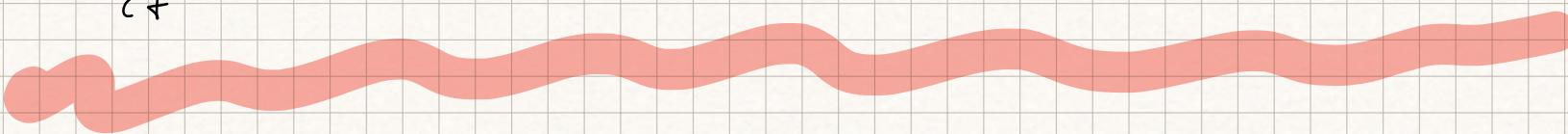
$$l_3 = \frac{8}{27}$$

$$\frac{20}{27} < 1 \rightarrow 0$$

$$\frac{40}{27} > 1 \rightarrow 1 \quad x = \frac{40}{27} - 1 = \frac{13}{27}$$

$$\frac{46}{27} < 1 \rightarrow 0$$

$$\frac{52}{27} > 1 \rightarrow 1$$



$$S = [1, 2, 3, 4, 6, 8, 9] \quad \text{interpolative code}$$

$l=1, r=4, low=1, high=9$

$$m = \left\lfloor \frac{1+4}{2} \right\rfloor = 4 \quad S_m = 4$$

$$Q = 1 + (4-1) = 4$$

$$S_{m-2} = 0 \quad m \quad \lceil \log_2(6-4+1) \rceil = 2$$

$$b = 5 - (4-4) = 6$$

(00)₂

$l=1, r=3, low=1, high=3$

$$m = 2 \quad S_m = 2$$

$$Q = 1 + (2-1) = 2$$

$$S_{m-2} = 0 \quad m \quad \lceil \log_2(2-2+1) \rceil = 0 \quad \text{no bit setted}$$

$$b = 3 - (3-2) = 2$$

$l=5, r=7, low=5, high=9$

$$m = 6 \quad S_m = 8$$

$$Q = 5 + (6-5) = 6$$

$$S_{m-2} = 2 \quad m \quad \lceil \log_2(8-6+1) \rceil = 2 \quad 10$$

$$b = 9 - (7-6) = 8$$

$L = 1110000111000111$

$H = \underbrace{001}_{0} \underbrace{101}_{1} \underbrace{01}_{2} \underbrace{110}_{3} \underbrace{110}_{4} \underbrace{0}_{5} \underbrace{10}_{6} \underbrace{0}_{7}$

$n = \# 1 \text{ in } H = 8$

$$l = \frac{L}{n} = \frac{16}{8} = 2$$

h	ℓ
0 0	1 1
0 1	1 0
1 0	0 0
1 0	0 1
1 0	1 1
1 0	0 0
1 0	0 1
1 1	1 1

$$h = \log_2 \# \mathcal{O} = 3$$

a a b a b c

C777 $\langle 0, 0, 0 \rangle \langle 1, 1, b \rangle \langle 2, 2, c \rangle$

C755 $\langle 0, 0 \rangle \langle 1, 1 \rangle \langle 0, b \rangle \langle 2, 2 \rangle \langle 0, c \rangle$

$\langle 0, \text{char} \rangle$ new char

$\langle d, \ell \rangle$ copy