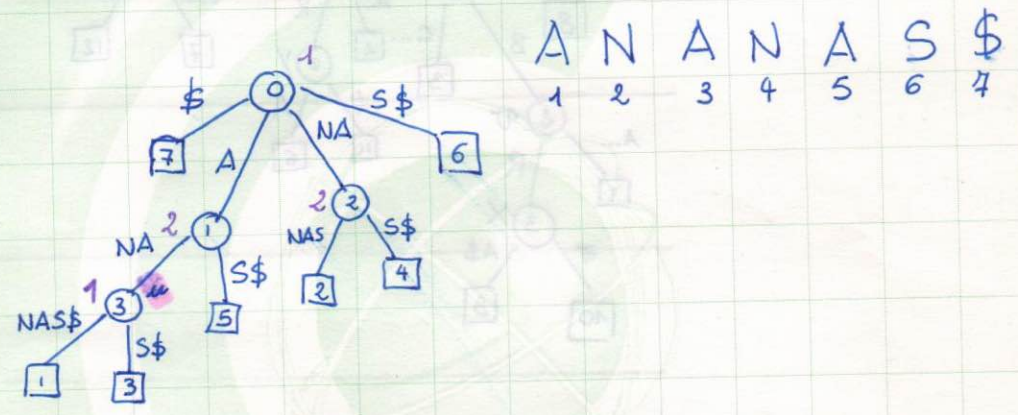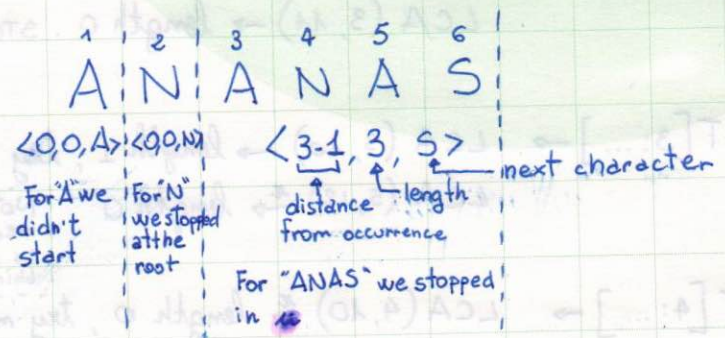# EXERCISES

## EXCERCISE 1

Let T = ANANAS $. Compute LZ - parsing using suffix tree.

1) Build the suffix tree (notice that it is typical to add the "$" at the end so no string is a prefix of all the others).

$$\begin{array}{ccccccc} A & N & A & N & A & S & \$ \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 \end{array}$$



2) Precompute minimum leaves for each node (purple color)

3) Percolate a path until the label is greater or equal than the position

$$\begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ A & N & A & N & A & S \end{array}$$

⟨0,0,A⟩ ⟨0,0,N⟩ ⟨3-1, 3, S⟩ → next character

For "A" we didn't start

For "N" we stopped at the root

⟨3-1, 3, S⟩
- 3-1 ↑ distance from occurrence
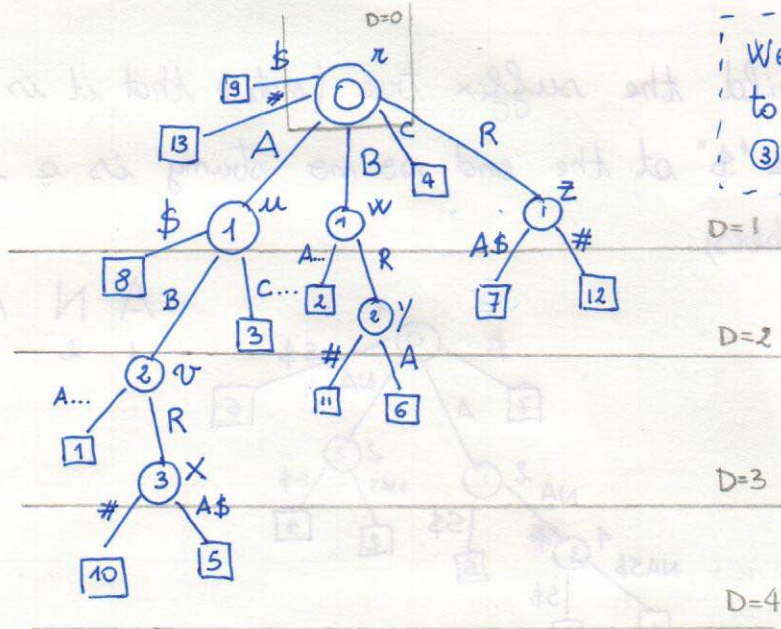- 3 ↑ length
- S ↑ the next character

For "ANAS" we stopped in

# EXERCISE 2

Let T = ABACABRA and let P = ABR. How does the 1-mismatch algorithm works on T and P, given that it is provided for free a data structure that solves LCA in $O(1)$ time. Only find 2 1-mismatch.

① Concatenate T and P : ABACABRA$ABR#
1 2 3 4 5 6 7 8 9 10 11 12 13

② Build suffix tree



We label internal nodes to ease notation in step ③

These Ds indicate the depth of nodes and are needed for the second part of the exercise (see next page)

③ The algorithm compares P with the first suffix of T

- $P \sim T[1: \dots]$ → LCA $(1, 10)$ has length 2, mismatch on (A, 3) and (R, 12) ✓

- $P \sim T[2: \dots]$ → LCA $(2, 10)$ → length 0, try next character
  LCA $(3, 11)$ → length 0. STOP because 2 mismatches

- $P \sim T[3: \dots]$ → LCA $(3, 10)$ → length 1, try next character
  LCA $(5, 12)$ →$x$→ length 0    Notice that since the lcp is 1 it means that the second character is a mismatch, so we jump to the third character

- $P \sim T[4: \dots]$ → LCA $(4, 10)$ →$x$→ length 0, try next character
  LCA $(5, 11)$ →$x$→ length 0 STOP because 2 mismatches

- $P \sim T[5: \dots]$ → LCA $(5, 10)$ →$x$→ length 3, no mismatch ✓

With respect to this exercise provide a data structure that allows LCA in constant time and compute LCA(10,6), block=4

① Perform Euler tour (without characters "$" and "#")

$$ET= r u 8 u \: v \: 1 \: v \: x \: 10 \: x \: 5 \: x \: v u \: 3 \: u \: r w 2 \: w y \: 11 y \: 6 y \: w x \: 4 x \: z 7 z \: 12 \: z x$$
$$D= 0 1 2 \: 1 \: 2 \: 3 \: 2 3 \: 4 \: 3 \: 4 3 \: 2 1 \: 2 1 \: 0 \: 1 2 \: 1 2 \: 3 \: 2 3 \: 2 \: 1 0 \: 1 0 \: 1 2 1 \: 2 \: 1 0$$

② Compute depth of each node ⎯

③ We recall that we have 3 pieces:
   1) D' blocked D with sparsification (powerbit)
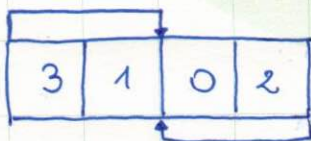   2) prefix-suffix minima of block of D
   3) Table of 0-1 s

> In this kind of exercise we are asked to draw only the data structure(s) we need to complete the exercise

④ Partition D into blocks of size 4 and highlight 10 and 6 and compute minima

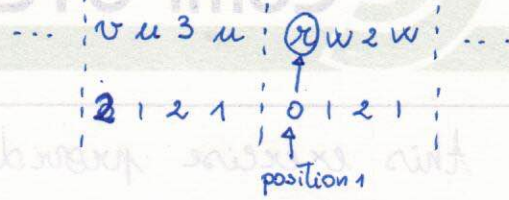| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|------|------|------|------|------|------|------|------|-----|
| 0121 | 2323 | 4343 | 2121 | 0121 | 2323 | 2101 | 0121 | 210 |
| 0 | 2 | 3 | 1 | 0 | 2 | 0 | 0 | 0 |

⑤ We notice that we need only sparsification, since we need to compute LCA of blocks 3, 4, 5, 6.

| 3 | 1 | 0 | 2 |
|---|---|---|---|

Notice that we take the biggest power of 2 smaller than the length

$$min \begin{cases} min \: D'[3 \: ; \: 3+2-1] = 1 \\ min \: D'[6-2^1+1 \: : \: 6] = 0 \end{cases} \longrightarrow 0$$

3

⑥ Find the relative position of the minimum in the block, namely 1 and find the corresponding node in the Euler tour:

$$\ldots \; \overline{v \; u \; 3 \; u} \; \fbox{$\overbrace{\textcircled{r} \; w \; s \; w}$} \; \ldots$$

$$\underline{2 \; 1 \; 2 \; 1} \quad \underline{0 \; 1 \; 2 \; 1}$$

$$\uparrow$$
position 1

Answer $r$.

Notice that the question ~~would~~ ~~this~~ ~~been~~ (was) to provide a data structure that ~~you~~ takes $n \cdot \log n$ the trivial answer would have been sparsification on the entire array.
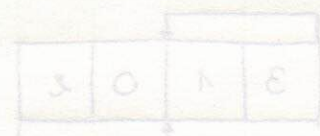
# EXERCISE 3

Provide the first 10 codewords of s.c. code, where s=1, c = 3.

Since $s + c = 4 = 2^2$, the number of bits needed is 2

| | | |
|---|---|---|
| 00 | STOPPER | |
| 01 | | |
| 10 | CONTINUERS | |
| 11 | | |

| | | | | |
|---|---|---|---|---|
| s { | 0 | 00 | | |
| cs { | 1 | 01 | 00 | |
| | 2 | 10 | 00 | |
| | 3 | 11 | 00 | |
| ccs { | 4 | 01 | 01 | 00 |
| | 5 | 01 | 10 | 00 |
| | 6 | 01 | 11 | 00 |
| | 7 | 10 | 01 | 00 |
| | 8 | 10 | 10 | 00 |
| | 9 | 10 | 11 | 00 |
| | 10 | 11 | 01 | 00 |
| | | 11 | 10 | 00 |
| | | 11 | 11 | 00 |

we stop here since we are required 10 cws

A variation of this exercise is: provide the configuration of number 15.

# of integers that can be written as $\boxed{s} \to 1$

# " " " " " " " " $\boxed{c}\ \boxed{s} \to 3 \cdot 1$

# " " " " " " " " $\boxed{c}\ \boxed{c}\ \boxed{s} \to 3 \cdot 3 \cdot 1$

# " " " " " " " " " $\boxed{c}\ \boxed{c}\ \boxed{c}\ \boxed{s} \to 3 \cdot 3 \cdot 3 \cdot 1$

It comes without saying that since $15 > 1 + 3 + 9 - 1$ the configuration is $c\ c\ c\ s$.

$\underset{\underset{\text{from 0}}{\underset{\uparrow}{\text{we start}}}}{-1}$

# EXERCISE 4

Let $S = 2\ 5\ 8\ 10\ 11\ 12\ 13$. Perform interpolative coding
on one level of recursion.

$\quad$ indices: $2_{(1)}\ 5_{(2)}\ 8_{(3)}\ 10_{(4)}\ 11_{(5)}\ 12_{(6)}\ 13_{(7)}$

① $\quad l = 1, \quad r = 7, \quad m = 7, \quad low = 2, \quad high = 13$

② $\quad$ Start encoding 10, because $m = \dfrac{l+r}{2} = \dfrac{7+1}{2} = 4$ and $S_m = 10$

③ $\quad$ Compute range $= [5, 10]$

> Do you recall?
>
> range $= [low+m-l, \ high-x+m]$

④ $\quad$ Encode 5 $\left(\substack{\text{offset of} \\ S_m \text{ in the} \\ \text{range}}\right)$ in $\lceil \log_2 \overset{\text{maximum gap}}{6} \rceil$ bits $\rightarrow$ 101

⑤ $\quad$ Recur on $2_{(1)}\ 5_{(2)}\ 8_{(3)}$ and on $11_{(5)}\ 12_{(6)}\ 13_{(7)}$

$\boxed{1}\ m=3\ l=1,\ low=2,\ r=3,\ high=9$ $\qquad$ $l=5,\ low=11,\ x=7,\ high=13,\ m=3$
$\qquad\qquad\qquad\qquad\qquad\qquad \underset{S_m-1}{\overset{4}{\uparrow}}$ $\qquad$ Since $13-11+1 = 3 = n$ no bits

$\boxed{2}\ m=2 \Rightarrow S_m = 5$ $\qquad\qquad\qquad\qquad$ are emitted for this part of $S$

$\boxed{3}\ \ $ range $= [3, 8]$

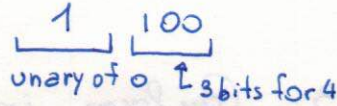$\boxed{4}\ \ $ Encode $5 - \text{\small(scribbled)} = 2$ in

$\qquad$ 3 bits $\rightarrow$ 010

# EXERCISE 5

(a) Rice code, where $x = 5$, $K = 3$

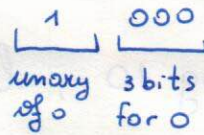① $q = \left\lfloor \dfrac{x-1}{2^k} \right\rfloor = 0$,     $r = 4$

② Concatenate $q$ in unary with binary in $k$ bits of $r$

$$\underbrace{1}_{\text{unary of } 0} \ \underbrace{100}_{3 \text{ bits for } 4}$$

(b) Compute $x = 1$, $K = 3$ with Rice code

① $q = \left\lfloor \dfrac{x-1}{2^k} \right\rfloor = 0$,     $r = 0$

② Concatenate as before     $\underbrace{1}_{\substack{\text{unary} \\ \text{of } 0}} \ \underbrace{000}_{\substack{3 \text{ bits} \\ \text{for } 0}}$

(c) Decode  001001 01 000 , given $K = 3$

① Split the sequence into block, knowing that we have a unary-binary sequence:

$$\underbrace{001}_{\text{unary}} \underbrace{001}_{\text{bin}} \Big| \underbrace{01}_{\text{unary}} \underbrace{000}_{\text{bin}}$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$q=2 \ , \ r=1 \qquad q=1 \ , \ r=0$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$1 + 2^3 \cdot 2 + 1 \qquad 1 + 2^3 \cdot 1$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$18 \qquad\qquad\qquad 9$$

Notice that the formula ~~as~~ depends on the ~~"~~ number $x-1$, ~~which~~ which may be divisible by $2^k$ or not.

# EXERCISES

## EXERCISE 1

We are given a text, where the frequencies are the following

| a | b | c | d | e | f | g |
|---|---|---|---|---|---|---|
| 0.1 | 0.1 | 0.11 | 0.1 | 0.1 | 0.21 | 0.28 |

. Build a classical Huffman code



Then, the exercise asks to build the canonical Huffman.

① Plot num array

| $\ell$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| num | 0 | 2 | 3 | 2 |

② Plot SYMB table

| $\ell$ \ symb | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 1 | | | | |
| 2 | f | g | | |
| 3 | c | d | e | |
| 4 | a | b | | |

③ Plot Fc array

| $\ell$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Fc | 2 | 2 | 1 | 0 |

Then the exercise asks to decompress the $\overset{\text{first three symbols of}}{\text{following}}$ sequence

$$C = \underbrace{0011}_{c}\underbrace{010}_{f}\ 010\ 10\ 110\ 111\ 00\ldots$$

$\ell = 1$ , $(0)_2 < $ fc$[1] = (2)_{10}$

$\ell = 2$ , $(00)_2 < $ fc$[2] = (2)_{10}$

$\ell = 3$ $(001)_2 = $ fc$[3] = (1)_{10}$ stop. 001 first codeword, which corresponds to the first character encoded by 3 bits: c

$\ell = 1$ , $(1)_2 \leq $ Fc$[1] = (2)_{10}$

$\ell = 2$ , $(10)_2 \overset{=}{=} $ fc$[2] = (2)_{10}$ stop. 01 first character encoded by 2 bits: f is the second codeword

Since the following to digits are 10 it comes without saying that the first three characters are cff.

# EXERCISE 2

Decode with arithmetic coding the ~~first three characters of the~~ following $c = 10011$, provided that $p(a) = \frac{1}{2}$ and $p(b) = p(c) = \frac{1}{4}$

① Diadic fraction : $1 \cdot \frac{1}{2} + 0 \cdot \frac{1}{4} + 0 \cdot \frac{1}{8} + 1 \cdot \frac{1}{16} + 1 \cdot \frac{1}{32} = \frac{19}{32}$
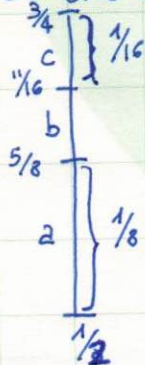
② Plot the intervals



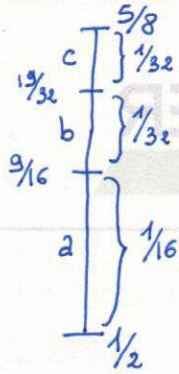Since $\frac{1}{2} < \frac{19}{32} < \frac{3}{4}$

we split the second interval (b)

> Always use fractions instead of floating points.
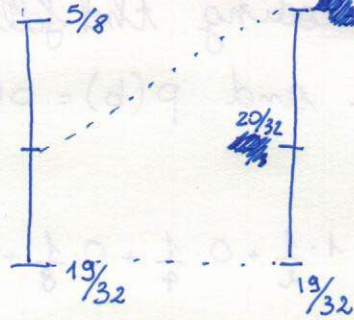> Always write letters sorted alphabetically

③ Plot the new intervals



Since $\frac{19}{32} < \frac{5}{8}$ we split the first interval (a)

3

④ Plot the new intervals



Since $\frac{19}{32}$ is represented as a border of the interval we should continue splitting the interval $[\frac{9}{16}, \frac{19}{32}]$. but we are required 3 symbols and we obtained: b a c

Observation: since the next intervals would have been we can say that the sequence is ba c ā
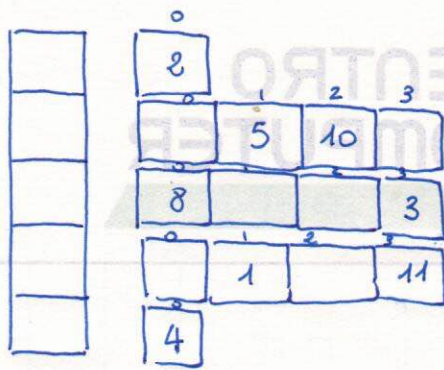
# EXERCISE 3

Compute the perfect hash for $S = \{1, 5, 11, 3, 10, 8, 2, 4\}$ provided that the first level $m = 5$.

① Pick at random a function $h_{ab}(x) = ax + b \mod p$, say $h_{2,1}(x) = 2x + 1 \pmod 5$

② Compute the number of collisions, suggestion: build graphically hashing with chaining

$$
\begin{array}{llll}
m_0 = 1 & 0 & \rightarrow & 2 \\
m_1 = 2 & 1 & \rightarrow & 5, 10 \\
m_2 = 2 & 2 & \rightarrow & 3, 8 \\
m_3 = 2 & 3 & \rightarrow & 1, 11 \\
m_4 = 1 & 4 & \rightarrow & 4 \\
\end{array}
$$

Is this choice ok? $\sum_{i=0}^{4} (m_i)^2 \stackrel{?}{<} 2m \iff 1+4+4+4+1 \stackrel{?}{<} 16$

$\iff 14 < 16$ ok

③ Let us build $(m_i)^2$ buckets $\forall i = 0, \ldots, 4$ and define 5 hash functions

$$h_0 = x \pmod 1$$

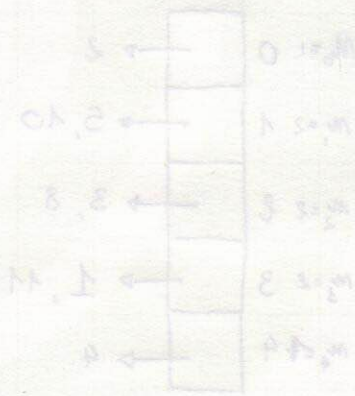$$\left.\begin{array}{l} \\ \end{array}\right\} \; h_1 = h_2 = h_3 = x \pmod 4$$

$$h_4 = x \pmod 1$$

When $m_i = 1$ we may choose $a = 1$, $b = 1$ and $m = 1$

## Perform search of Key 7.

① $7 \cdot 2 + 1 \equiv 15 \equiv 0 \pmod 5 \;\to\; 1° \text{ bucket}$

② $7 \equiv 0 \pmod 1 \;\to\; 1° \text{ bucket}$
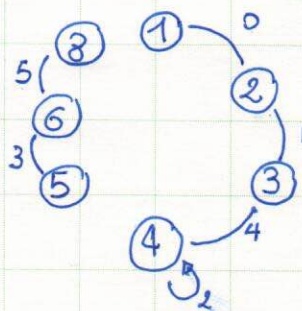
③ $\boxed{2}$ does not contain 7, hence 7 is not present.

# EXERCISE 4

Given a set of strings $S = \{aa, ab, bb, bc, ca, db\}$, provided that the rank of the symbols is $a \to 1$, $b \to 2$, $c \to 3$, $d \to 4$ and $h_1(xy) = r(x) \cdot r(y) \pmod{m}$ and $h_2(xy) = r(x) + r(y) \pmod{m}$, build a MOPHF (Minimal Ordered Perfect Hash Function), on $m = 11$.

① Compute $h_1$, $h_2$ and $h$ for all the keys

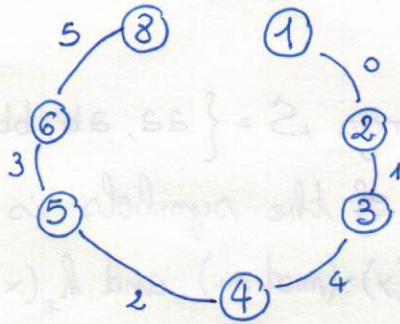|     | $h_1$ | $h_2$ | $h$ |
|-----|-------|-------|-----|
| aa  | 1     | 2     | 0   |
| ab  | 2     | 3     | 1   |
| bb  | 4     | 4     | 2   |
| bc  | 6     | 5     | 3   |
| ca  | 3     | 4     | 4   |
| db  | 8     | 6     | 5   |

② Build the graph



> Only draw the nodes that are values of $h_1$ and $h_2$

Since the graph is not acyclic it is not possible to build a MOPHF.

The exercise asks to modify the graph to allow the creation of a MOPHF.
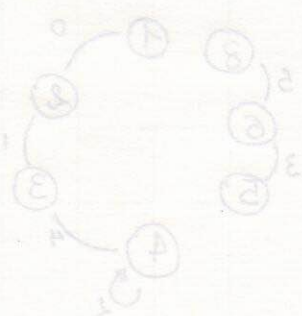
① Eliminate self loop, we choose



② Define by cases the function $g$

| $g$ | $x$ |
|-----|-----|
| 0 | 1 |
| 0 | 2 |
| 1 | 3 |
| 3 | 4 |
| 5 | 5 |
| 4 | 6 |
| / | 7 |
| 1 | 8 |

Notice that this method works only on keys in set $S$.

Let us take $ac$, $h_1(ac) = 3$, $h_2(ac) = 4 \Rightarrow g(h_1(ac)) = 1$, $g(h_2(ac)) = 3$, hence the algorithm would assign to "$ac$" rank 4, which is wrong.

# EXERCISE 5

Assume we want to build a Bloom filter of an set of $n = 2^{16}$ keys building a binary array of $m = 2^{20}$ bits.

(a)_Compute the optimal number of hash functions and the optimal error.

(b)_Compute the size of the Bloom filter for $n = 2^{16}$ to guarantee $\varepsilon = 2^{-20}$ by using an optimal number of hash functions.

(a) $\quad K_{opt} = \dfrac{m}{n} \cdot \ln 2 = 2^4 \cdot \ln 2$

$\quad\quad \varepsilon_{opt} = (0.6185)^{16}$

(b) $\quad$ Set $2^{-20} = (0.6185)^{\frac{m}{2^{16}}}$ and solve the equation