

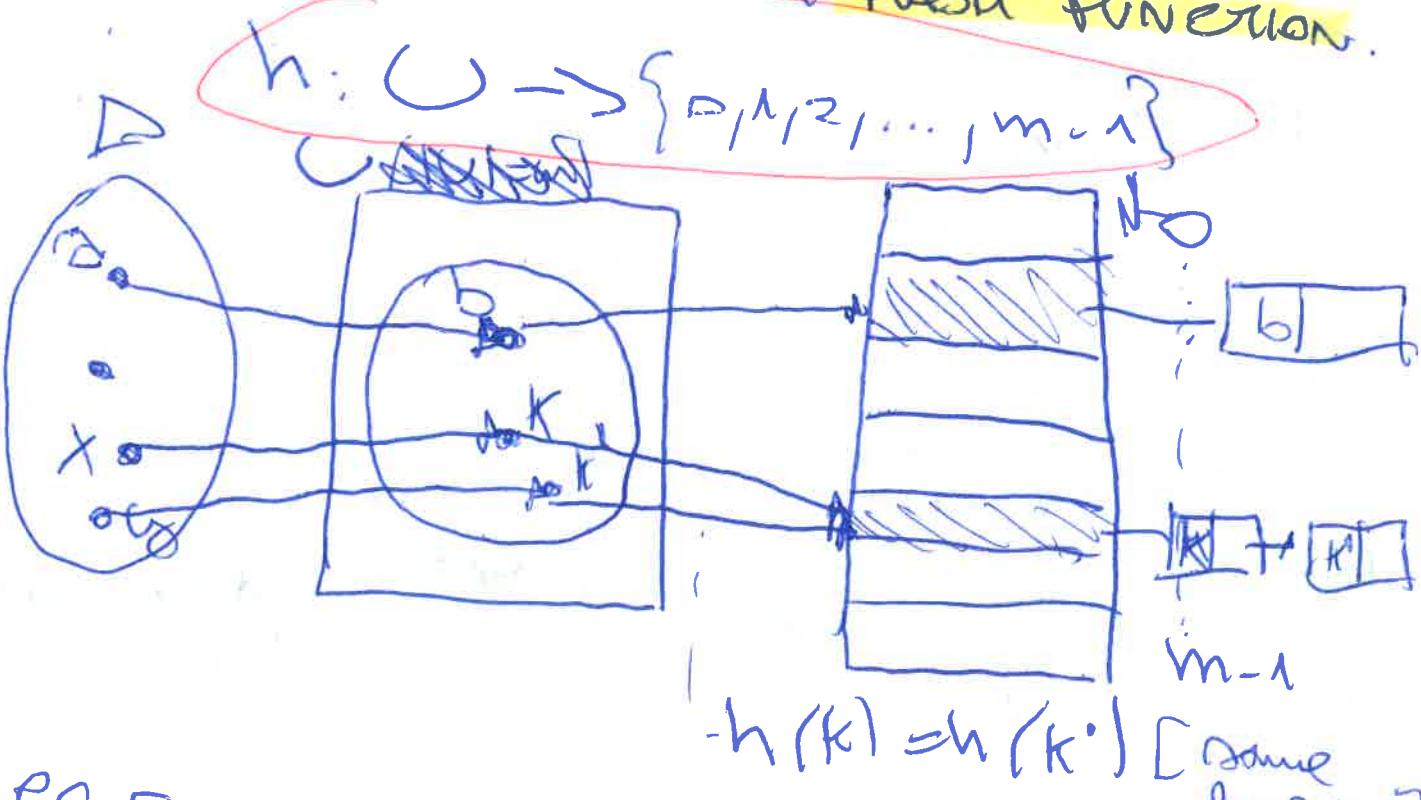
HASHING - HASH TABLES

(8)

HASHING WITH CHAINING:

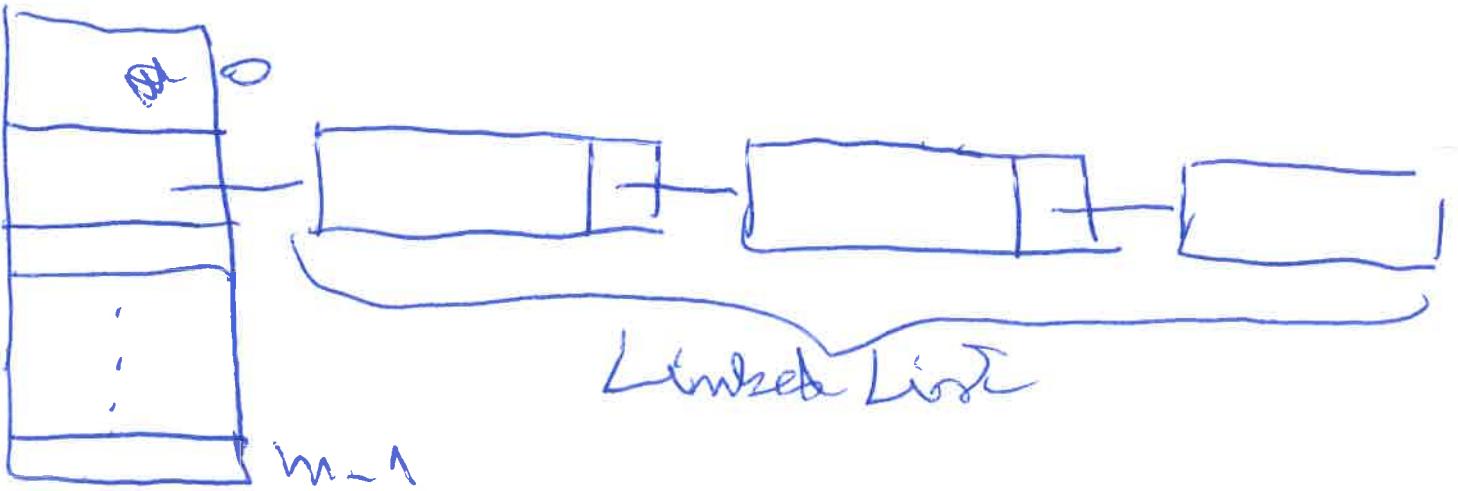
The hash table is an array of size T_m , whose entries are either NULL or they point to lists of dictionary items.

The mapping of items to array entries is implemented via a **HASH FUNCTION**.



PROPERTIES OF HASH FUNCTION h :

- h is a SIMPLE, UNIFORM hash,
(i.e. it can map any value to any possible cell with the same probability)
UNIFORM PROBABILITY.



$$E[\text{list length}] = \frac{N}{m} = \frac{\# \text{keys}}{\text{TABLESIZE}} = d = \frac{\text{load factor}}{\text{missing values}}$$

Generally, we have hash functions of the following form:

$$h(k) = k \bmod p$$

Do such functions always satisfy the UNIFORMITY property? Not

REQUIREMENTS for a HASH FUNCTION:

- FAMILY of HASH FUNCTIONS.

$$h \in \mathcal{H}: m_1, m_2, \dots, m_n \rightarrow m$$

FAMILY of SIMPLE, UNIFORM HASH FUNCTIONS

- In a family of hash functions, a hash function can map arbitrarily all keys to all possible values.

SPACE for the REPRESENTATION (in Bits)

$\forall h \in \mathcal{H}$, we need $\geq \log_2 m$ Bits,

$\cup \log_2 m$ Bits

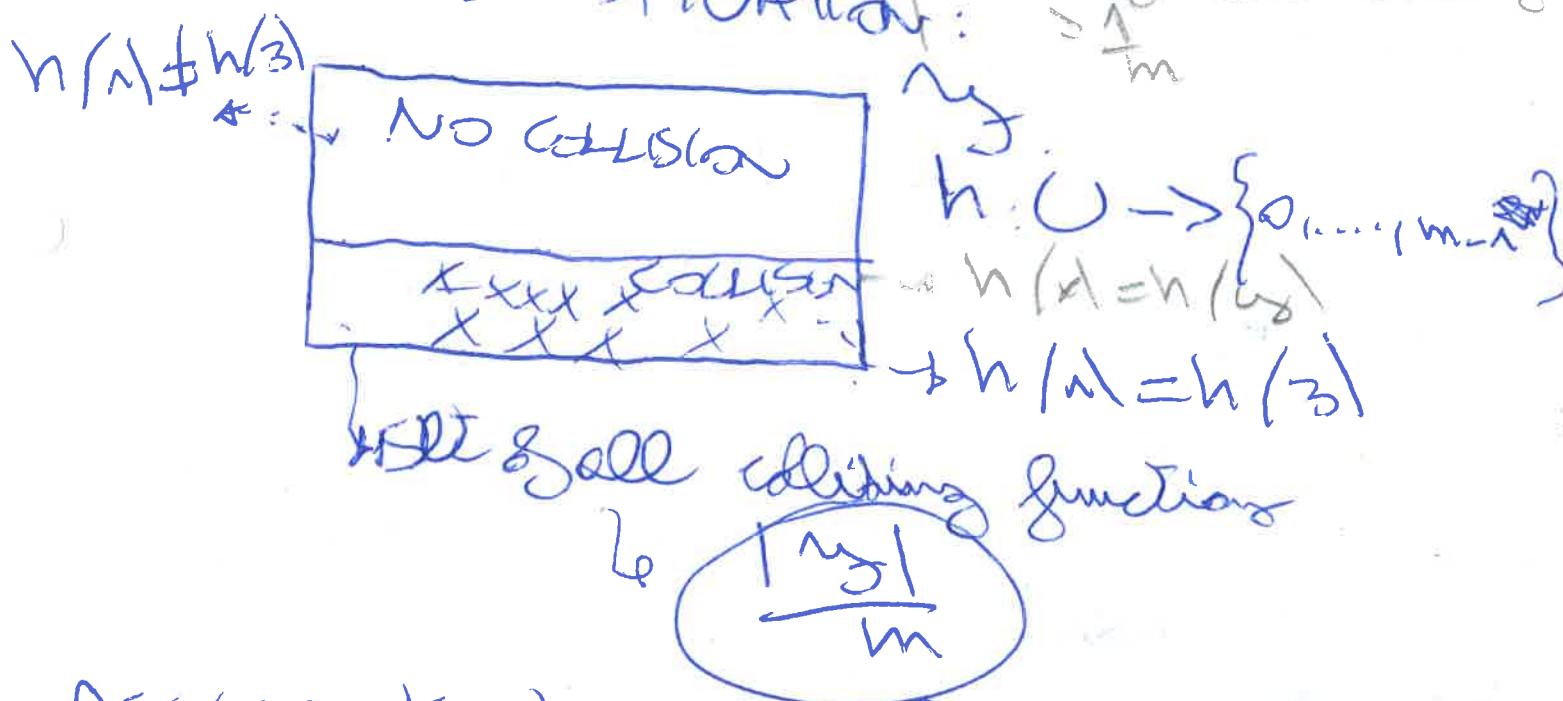
My PRACTICAL REPRESENTATION of every single HASH FUNCTION.

UNIVERSAL HASHING:

A family \mathcal{H} of hash function is a UNIVERSAL HASHING iff..

$$\forall x_1, x_2 \in U: [\#h: h(x_1) = h(x_2)] \leq \frac{1}{m}$$

MISERIALIZED SITUATION: (i.e. P. of collision occurring)



ASSUME that one such family \mathcal{H} exists.

\Rightarrow What is the P. to pick a hash function (at random) that would make keys collide?

(i.e. one ends up in the "collision part")

~~Probability of collisions~~

Pick $h \in \mathcal{Y} = \{h(x) = h(y)\}$

$$= \frac{\#\{x\}}{m} \text{ colliding hash functions} = \frac{1}{m} = P_{\text{collide}}$$

= $\frac{1}{m}$ hash function picked at random makes ~~any~~ collision

DEFINITION OF UNIVERSAL HASHING:

In a UNIVERSAL HASHING FAMILY, a randomly-chosen hash function h from the set \mathcal{Y} has P . to make distinct keys x and y to collide no more than $\frac{1}{m}$.

LENGTH OF CHAINING LISTS IN UNIVERSAL HASHING:

Let $T[0:m-1]$ be a hash table with chaining, where the hash function is picked randomly from a universal class \mathcal{Y} .
 \Rightarrow the expected length of the chaining lists is still no more than $d + \alpha$, where
 [i.e. a hash function taken from the family \mathcal{Y} still has load factor $\alpha = \frac{n}{m}$]

$$\alpha = \frac{n}{m}$$

$$d = \lceil \frac{n}{m} \rceil$$

PROOF:

→ The expected load factor depends on the choice of h from \mathcal{H} where \mathcal{H} is a FAMILY OF UNIVERSAL HASH FUNCTIONS.

→ We hence define:

$$I_{X_{k,y}} = \begin{cases} 1 & \text{if } h(x) = h(c_y) \text{ (keys collide)} \\ 0 & \text{Otherwise} \end{cases}$$

V $X_{k,y}$ is
a indicator
RV.

Where $I_{X_{k,y}}$ is an INDICATOR R.V.

~~We want to determine the # keys where we have a collision, i.e.~~

~~# $y \in D \mid h(x) = h(c_y)$~~

We are interested in determining the length of the chain for finding the Load Factor.

i.e.: $\sum_{y \in D} I_{X_{k,y}}$ = Length of chain.

where * $y \in D : h(x) = h(c_y)$

$$E\left\{\sum_{y \in D} I_{X_{k,y}}\right\} \leq \sum_{y \in D} E\{I_{X_{k,y}}\}$$

$$E\{I_{X_{k,y}}\} = P\{I_{X_{k,y}}=1\} + P\{I_{X_{k,y}}=0\} = 0$$

$$\Rightarrow E\{I_{X_{k,y}}\} = P\{I_{X_{k,y}}\} \cdot 1$$

$$\exists \text{ such } \sum_{y \in D} E\{I_{X,y}\} = \sum_{y \in D} P\{I_{X,y}=1\}$$

$$= \sum_{\substack{y \in D \\ (y \neq x)}} P\{h(x) = h(y)\} = \frac{n-1}{m}$$

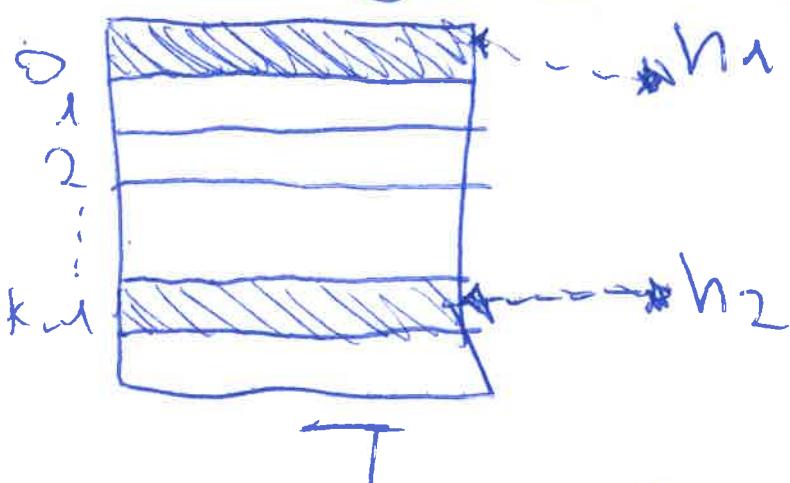
because g_x
 including x
 $\frac{n-1}{m} = \frac{1}{2}$

\Rightarrow We have shown that we STILL have the choosing hashing behaviour.
 (or A.V.A.)

TIME COMPLEXITY → Basic FUNCTION From UNIFORM FUNCTIONS' FAMILY

- We know:
 - Avg. length of a LIST = $O(\frac{n}{m})$ = $\frac{n}{m}$ most-use
 - MAX. length of a LIST = $O(\frac{\log n}{\log \log n})$
- Even if not constant time, we still have very little time.

* The Power of 2 Choices:



WEBSITE IDEA: Compute hash value by 2 hash functions; then ~~choose the lowest~~

the hashed value into the chain first
has the last LOAD FACTOR $\lambda = \frac{n}{m}$

SPACE = $O(\log \log n) = \Theta(m)$

(previously $\log n \rightarrow \log \log n$)

reduced by one order of magnitude



TIME = $O\left(\frac{\log n}{\log \log n}\right)$

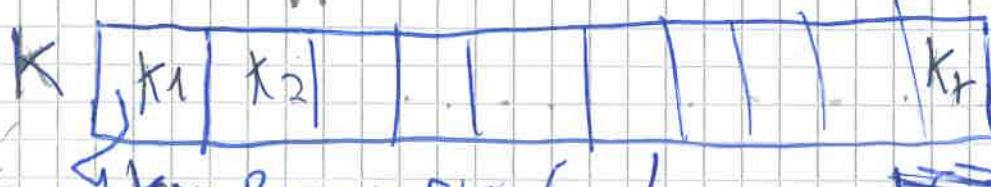
NO READS!

NB: Operating Systems use this kind of
format to reduce space occupancy ~~esp.~~
with Hash tables.

(~~No~~ PROOF!)

UNIVERSAL HASH FUNCTIONS DESIGN

We are interested in designing a hash FUNCTION that is UNIVERSAL.



SAME $\rightarrow h = \log_2 m \text{ bits}$



(# elements to work into T)

Bits to represent the hash UNIVERSE
"log2 U", where $U = \frac{\log_2 U}{\log_2 m}$

FAMILY OF HASH FUNCTIONS (with parameter m)

A family of hash functions with parameter m is defined as:

$$H = \{h_a : U \rightarrow \{0, 1, \dots, m-1\}\}$$

(i.e.: for every integer parameter a , you have the following hash function):

$$h_a(k) = \sum_{i=1}^r a_i k_i \bmod m$$

UNIVERSAL
HASH FUNCTION

STRUCTURE VISUALIZED

POINT



"log2(U)"

Key aspect: $k = [k_1 | k_2 | \dots | k_r]$
into $\log_2 m$ bits

Key k is split into r parts

$\log_2 m$ bits

parts of $\log_2 m$

UNIVERSAL HASH FUNCTION'S DEFINITION.

~~h₂~~ is a UNIVERSAL HASH FUNCTION.

$$\{ \# h_2 : h_2(x = h_2(y)) \leq \frac{1}{m} \}$$

(The hash functions where we have a collision is $\sum \frac{1}{m}$)

for a random $h \in H$
chosen from \mathbb{Z}_m^m , the
 $P[h(x) = h(y)] \geq \frac{1}{m}$

PROOF OF THE EXISTENCE OF UNIVERSAL HASH FUNCTIONS

We want to show that the inequality is indeed satisfied: Two distinct keys

$m \rightarrow$ PRIME NUMBER.

For $x_0 \neq y_0$ we about a counterexample (omit the cases where this occurs)

$$h_2(x) = h_2(y) \leftrightarrow \sum_{i=0}^{t-1} a_i x_i \equiv \sum_{i=0}^{t-1} a_i y_i \pmod{m}$$

$$\leftrightarrow 2a_0 x_0 + \sum_{i=1}^{t-1} a_i x_i \equiv 2a_0 y_0 + \sum_{i=1}^{t-1} a_i y_i \pmod{m}$$

$$\leftrightarrow 2a_0(x_0 - y_0) \equiv \sum_{i=1}^{t-1} a_i(g_i - x_i) \pmod{m}$$

{ Because $x_0 \neq y_0$ and $y_0 \neq 0$ }

$$\Rightarrow x_0 - y_0 \neq 0$$

$$2a_0 \equiv \frac{1}{x_0 - y_0} \sum_{i=1}^{t-1} a_i (y_i - x_i) \pmod{m}$$

\Rightarrow The last equation obtained for z_0 shows that:
 whatever is the choice for $[z_1, z_2, \dots, z_{r-1}]$
 there exists only one choice for z_0

* that causes x and y to COLLIDE.
 $(x, y, h(x) = h(y))$.

\Rightarrow What are the # choices for $[z_1, z_2, \dots, z_{r-1}]$
 that cause x and y to collide?

$[z_1, z_2, \dots, z_{r-1}] \Rightarrow$ We have $r-1$
 choices of m !

Choices = $\underbrace{m \times m \times \dots \times m}_{r-1 \text{ times}}$

$$\text{for } \star_2 = m^{r-1} = \frac{m^r}{m} = \frac{|y|}{m} = \frac{1}{m}$$

~~O.E.W.~~

~~for all $i \in \{1, 2, \dots, r\}$~~

$$\Rightarrow \{ \# h_i : h_i(x) = h_i(y) \} \leq \frac{|y|}{m} \quad \text{Q.E.D}$$

(As we have as many hash functions
~~as parameter~~ $\textcircled{1}$ as parameters $\textcircled{2}$)

\Rightarrow i.e. the definition ~~is~~ holds out
 at the beginning ~~is~~ holds!

TWO-UNIVERSAL HASH FUNCTIONS.

We now check whether there are other ways of constructing hash function, hash w.r.t parameter 2.

Consider $|U| = 2^h$ and $m = 2^l \leq |U|$
($\& p: a power of 2$)
 $2 = \text{an integer} < 2^h$

$$h_0(k) = \underbrace{(2k \bmod 2^h)}_{\text{CO-PRIME}} \text{ DIV } 2^{h-l}$$

is TWO-UNIVERSAL HASH FUNCTION.

where $2^h \geq 2^l$

No proof of correctness, though we have..

$$P\{h_0(x) = h_0(y)\} \leq \frac{1}{2^{l-n}} = \frac{2}{m}$$

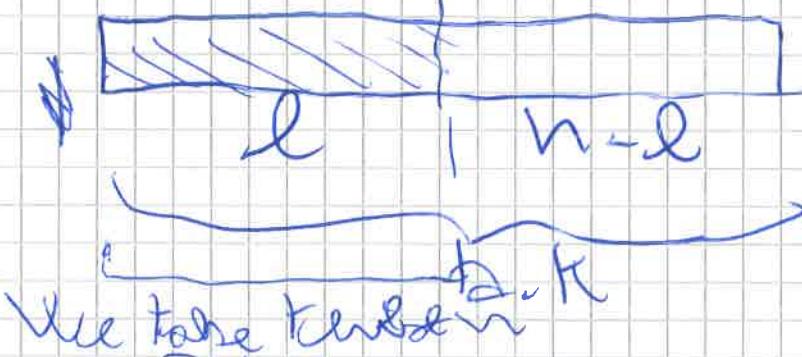
RECALL

N.B. We can also shuffle keys & use a mask to take the relevant bits,

or draw \Rightarrow

$P\{ \text{a collision occurs in a } \text{UNIVERSAL HASH FUNCTION} \} = \frac{1}{m}$

$P\{ \text{a collision occurs in a } \text{THE-UNIVERSAL HASH FUNCTION} \} = \frac{2}{m}$



PERFECT HASHING, MINIMAL ORDERED

SO FAR, Hashing with Chaining allows for storage in $\Theta(m)$ time on AVERAGE, where ~~where~~ $\ell = \Theta(m)$ if $N \approx m$.

(\rightarrow constant time on Average!)

In the Worst Case, we would have LINEAR

TIME access

FEATURES &

FIFO

\Rightarrow **PERFECT HASH**: We want to achieve $O(1)$ in the TIME access, with NO waste of space ($m = n$) still!

(Optimal) Linear Space \Rightarrow Constant Search Time \Rightarrow Access

PERFECT HASH FUNCTION - DEFINITION:

A hash function is said to be perfect with respect to a dictionary D iff:

~~# was~~ $h(x) \neq h(y) \Leftrightarrow$ ~~No~~ Collision among elements of D

ORDER-PRESERVING HASH FUNCTION

A (minimal) perfect hash function is said to be order-preserving if: $k_i < k_j \Rightarrow h(k_i) < h(k_j)$

$$\forall k_i, k_j \in S: h(k_i) < h(k_j) \quad (1)$$

$\Rightarrow h(k)$ return Rank of k in the ordered dictionary S . [Only in STATIC SCENARIO case]

MINIMAL HASH FUNCTION

A hash function is minimal if its size Hash Table size is $m = n$.

[Smallest possible size of Hash TABLE].

~~LEMMAS TO DESIGN A HASH TABLE~~

~~If T consists of q leaves and $m = q^2$ then a UNIVERSAL HASH FUNCTION.~~

DESIGN OF A SIMPLE & PERFECT HASH TABLE:

If ordering and minimality is NOT required, then the design of a (static) ~~perfect~~ hash function is simpler.

Use a TWO-LEVEL HASHING SCHEME with universal hashing functions at both levels.

TWO- LEVEL HASHING SCHEMES

To design a SIMPLE & PERFECT hash table, we make use of 2 LEVELS of hashing.

- * **FIRST- LEVEL Hashing:** Analogously to hashing with chaining, n keys^{NOTE} are hashed into m slots via a UNIVERSAL HASH FUNCTION h .
- * **SECOND- LEVEL Hashing:** Every entry of the first hash table points to a SECONDARY hash table T_2 through T_2^j via an own UNIVERSAL HASH FUNCTION h_2^j .

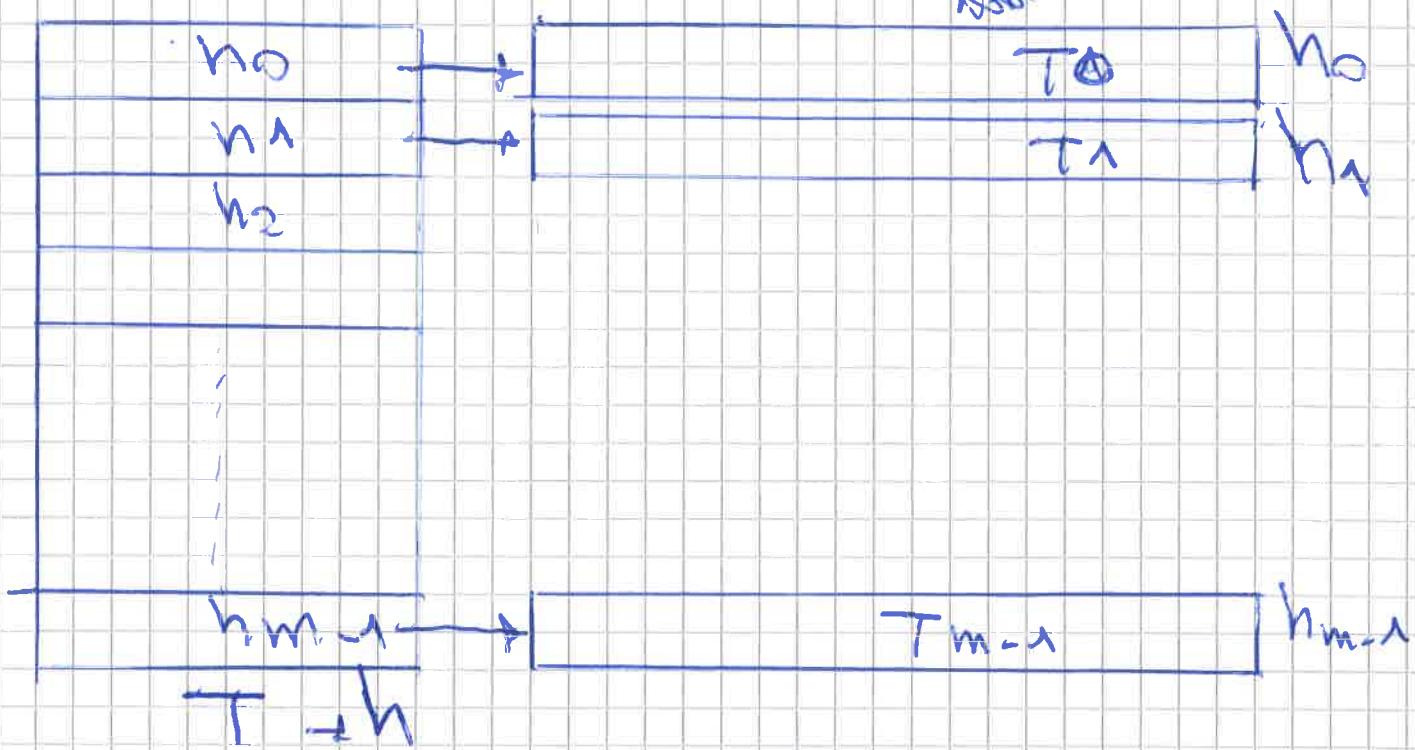
→ Searching for a key will require 2 TABLE ACCESSSES ~~→~~ 2 usage of 2 HASH FUNCTIONS.

$$2 \cdot n \text{ to } T_1 \{ O(n) \}$$

~~Worst case~~

$$2 \cdot h_i \text{ to } T_2^j \{ O(1) \}$$

VISUALIZATION:



We would like to ensure:

I) Space taken by T and all T_j 's suitable is $O(n)$.

II) The hash functions used to map T to T_j 's are PEFER.

LEMMA:

$$q = n$$

If we store q keys in a hash TABLE of size $m = q^2$, ~~using~~ a UNIVERSAL HASH FUNCTION, ~~then~~ then we have that:

$$\mathbb{E}\{\# \text{collisions}\} < \frac{1}{2} \quad \text{and}$$
$$\mathbb{E}\{\# \text{collisions}\} \leq \frac{q^2}{2m}$$

PROOF:

In a set of \textcircled{q} elements (Keys), we have

$$\mathbb{E}\{\# \text{collisions}\} = \binom{q}{2} < \frac{q^2}{2} \quad \text{because:}$$

$$\binom{q}{2} = \frac{q!}{(q-2)! \cdot 2!} = \frac{q \cdot (q-1)}{(q-2)! \cdot 2!} = \frac{q(q-1)}{2!}$$

$$\Rightarrow \frac{q \cdot (q-1)}{2} < \frac{q^2}{2} \quad \text{so why?}$$

If we pick \textcircled{n} function from a UNIVERSAL class of hash functions, then we have that each pair collides with $P = \frac{1}{m}$

$$\mathbb{E}\{\# \text{collisions}\} = \binom{q}{2} \cdot \frac{1}{m} < \frac{q^2}{2m}$$

And since we have taken $m = q^2$

$$\Rightarrow E\{\text{#collisions}\} = \binom{q}{2} \cdot \frac{1}{m} < \frac{q^2}{2}$$

$$\Rightarrow E\{\text{#collisions}\} < \frac{1}{2} \quad \text{Q.E.D.}$$

~~Upper Bound~~

\Rightarrow The probability to have at least one collision can be upper bounded via the Markov's INEQUALITY:

$$P\{X \geq t \cdot E\{x\}\} \leq \frac{1}{t},$$

where in our case $t=2$ and

$X = \#\text{collisions}.$

And also, we have shown that

$$P\{\text{collision}\} \leq \frac{1}{2}$$

$$\Rightarrow P\{\text{NOT P wave collision}\} \geq \frac{1}{2}$$

\Rightarrow We use this result to get

$$\text{I) } m \geq T_2 \approx n^2$$

SIZE of primary hash table

$$m \geq (n)^2$$

SIZE of secondary hash table
square of
#log2(m)
 $T(\sqrt{s})$

II) $m = n$ for T_1 to ensure that every hash function is perfect.

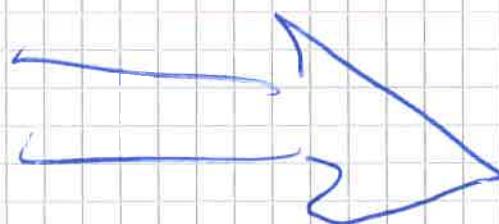
$\Rightarrow O(n)$ space. (very poor)

THEOREM - BOUNDING THE SPACE OF PRIMARY & SECONDARY HASH TABLES.

If we store m keys in a hash table of size $m = n$ using a UNIVERSAL HASH FUNCTION, then the expected size of all sub-tables T_j is $\leq 2n$.

PROOF. (SECONDARY HASH TABLES)

We operate according to the following procedure when designing a PERFECT HASH TABLE.



First-Level Hash Table

- (1) Distribute keys of T into T into n slots by a UNIVERSAL HASH F.
- (2) Let $\sum_{i=0}^{n-1} (n_i)^2 \leq 2n$, then

Keys leading to slot i Go to (1) and choose another hash function ~~for slot i~~

else (2) does not hold [we have too many collisions]

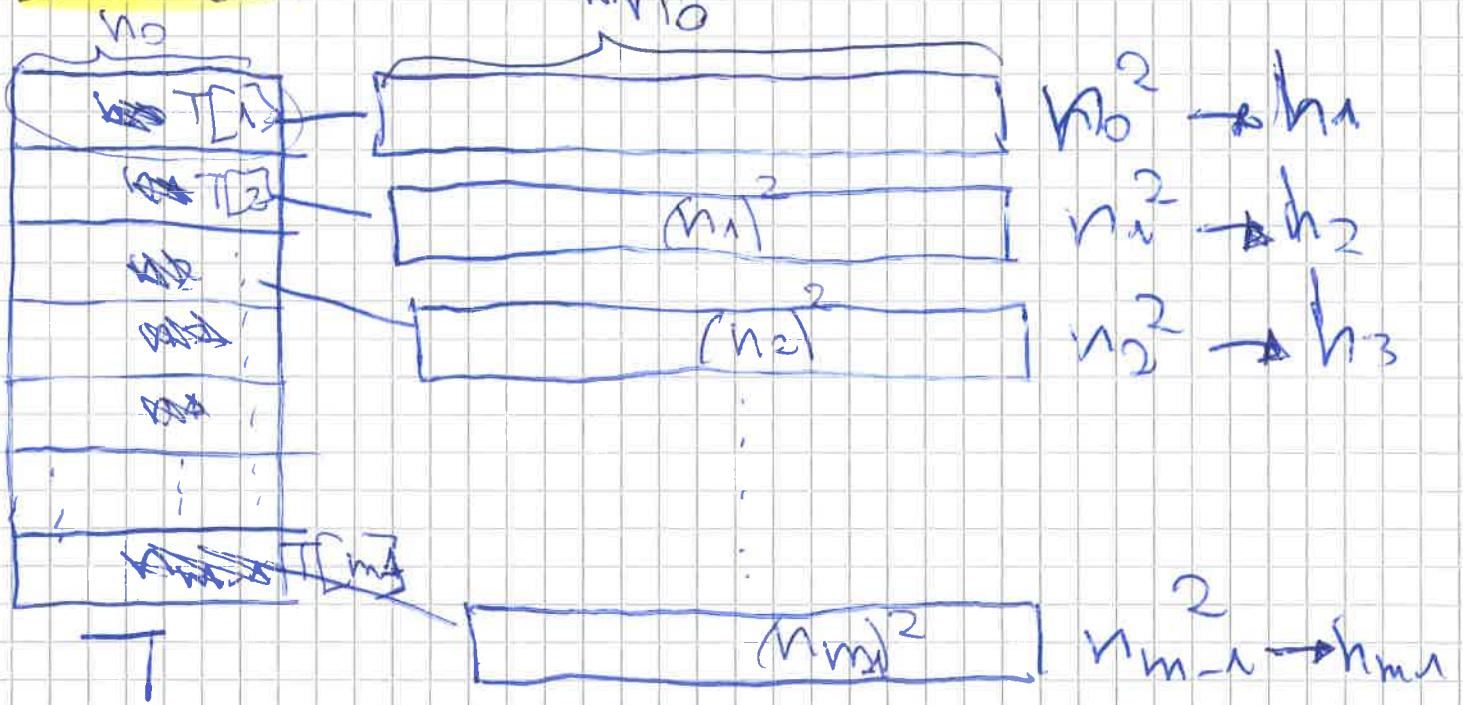
$$\sum_{i=0}^{n-1} (n_i)^2 \leq 2n$$

where $n_i = \# \text{keys leading to slot } i$

~~(1) & ~~choose another hash function~~~~

Primary Hash Table

Situation:



The different slots pick their own hash function in an INDEPENDENT manner.

⇒ Let's hence estimate the # collisions in the different slots

$$\sum_{i=0}^{n-1} (n_i)^2 = \sum_{i=0}^{n-1} (n_i) + 2 \cdot \frac{(n_i)}{2}$$

$$E\left\{\sum_{i=0}^{m-1} (n_i)^2\right\} = \sum_{i=0}^{m-1} (n_i) + 2E\left\{n_i\right\} = N^2$$

Because:

$$\begin{aligned} \sqrt{N^2} &= \sqrt{N + 2 \cdot \frac{N}{2}} = \sqrt{N + 2 \cdot \frac{N!}{(N-1)!}} \\ &= \sqrt{N + N \cdot \sqrt{N-1}} \\ &= \sqrt{N + N^2 - N} \end{aligned}$$

$$E\left\{2\sum_{i=0}^{m-1} (n_i)^2\right\} = E\left\{\sum_{i=0}^{m-1} n_i\right\} + 2E\left\{\sum_{i=0}^{m-1} \frac{n_i}{2}\right\}$$

~~total # atoms~~
~~masses and forces~~
~~masses and forces~~ SECONDARY WASH TUBES

$$= N + 2 \cdot E\left\{\sum_{i=0}^{m-1} \frac{(n_i)}{2}\right\}$$

What is $\frac{(n_i)}{2}$?

\Rightarrow It's the total # collisions

Consider dry primary LEVEL WASH

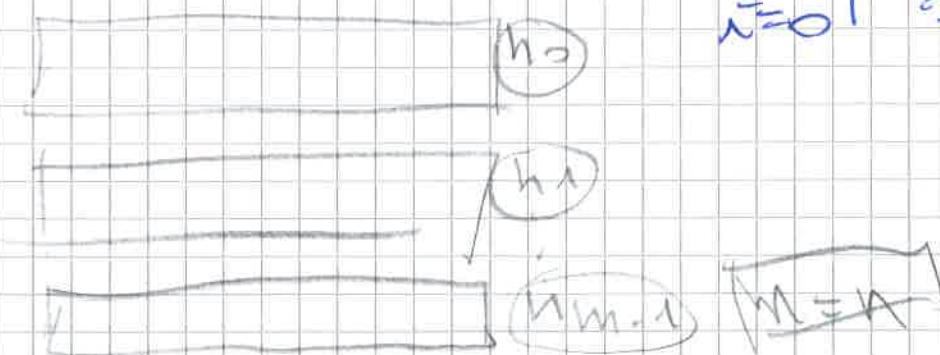
FUNCTION (then consider ~~every~~ every different set)

$$\sum_{i=0}^{m-1} \frac{(n_i)}{2}$$

$$\sum_{i=0}^{m-1} \frac{(n_i)}{2}$$

N # collisions in:

$$\Rightarrow \sum_{i=0}^{m-1} \frac{(n_i)}{2} = \binom{N}{2} \frac{1}{m}$$



$$\begin{aligned}
 \Rightarrow \text{DEG} \left\{ \sum_{i=0}^{m-1} \frac{1}{(n-i)} \right\} &= n+2 \cdot \binom{n}{2} \cdot \frac{1}{\frac{1}{n(n-1)}} \\
 m &= n \\
 &= n+2 \cdot \frac{n!}{(n-2)! \cdot 2!} \cdot \frac{1}{\frac{1}{n(n-1)}} \\
 &= n+2 \cdot \frac{n(n-1)}{2!} \cdot \frac{1}{\frac{1}{n(n-1)}} \\
 &= 2n-1
 \end{aligned}$$

\Rightarrow # collisions $\approx 2n-1 \approx 2n$

\Rightarrow We have hence shown that two-level hashing \Rightarrow to bound by $2n$.

On average, two iterations are sufficient to find a good hash function.

Also! We have shown that space $\approx 2n$

\Rightarrow space = $O(n)$

CREATING A MINIMAL, ORDERED PERFECT HASH FUNCTION:

Recall:

MINIMAL HASH FUNCTION:

$m = n$ (i.e. no wasted space)

ORDERED HASH FUNCTION:

by $x \leq y$:

$$\boxed{h(x) < h(y)}$$

$$\Leftrightarrow h = \text{Ranks}(K)$$

We want a function that, when given a key, it returns the rank [i.e. the index]

	h_1	h_2	rank
t_1	1 5	6 *	0
t_2	7	2	1
t_3	5	7	2
t_4	4	6	3
t_5	1	10	4
t_6	0	1	5
t_7	8	11	6
t_8	11	9	7
t_9	5	3	8

h_1 and h_2 are two **INDEPENDENT** hash functions (either given or can be arbitrarily chosen)

↳ Alternatively, we can just generate numbers for the table (randomly)

~~extra~~ Let's now consider also:

- \circ , a function that maps integers in range $\{0, \dots, m-1\}$ to $\{0, \dots, n-1\}$

~~What~~

We consider $m' = c \cdot n$, where

$$\begin{array}{c} m' \\ \leftarrow \rightarrow n \end{array}$$

two trials n_1
 $c = 3$

\Rightarrow \circ won't perform a **PERFECT MAPPING**, as $m' > n$.

We want to find Δt :

$$n(t) = \circ(h_1(t)) + \circ(h_2(t)) \bmod n$$

To do this, we make an **analog**.

Constructor ^{INDIRECTO} **CREATE** based on table,

where:

- ↳ Nodes are values of n_1 and n_2
- ↳ Edges connect the different nodes, and they are labeled by \circ (the resulting sum value).

② Find values for the ~~edges~~

$$\circ(h_1(n_1(6))) \text{ and}$$

$$\circ(h_2(n_2(6)))$$

$\Rightarrow x: 0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11$

$\circ: 0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11$

EXAMPLE.

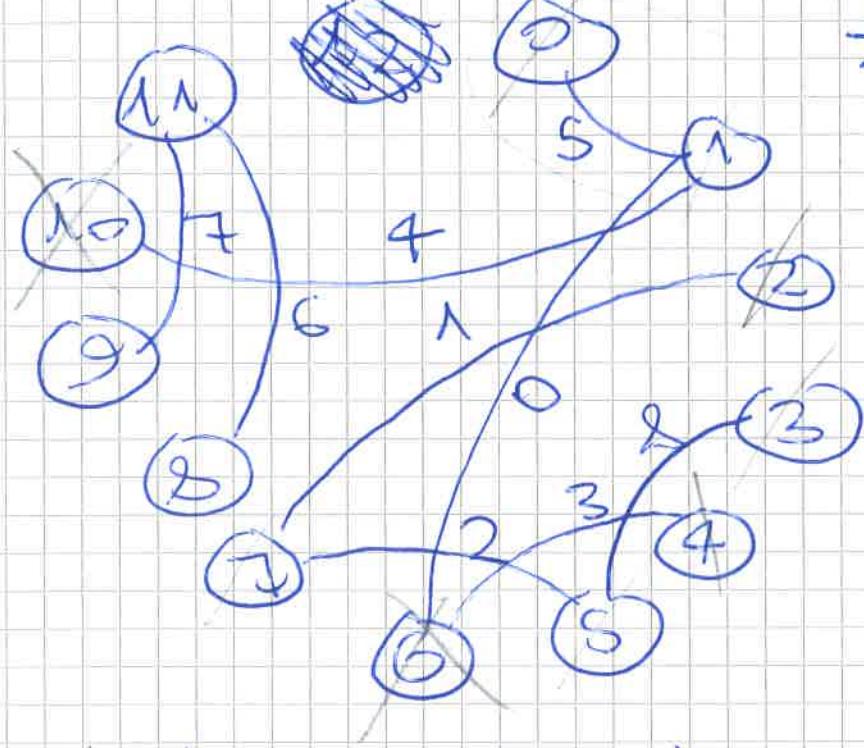
$C_0, m' - n \}$

$$m' = 12$$

(actually $m' = 13$)

$$n = 2.9$$

$$\# \text{NODES} = m'$$



$$h(\ell) = g(h_1(\ell)) + g(h_2(\ell)) \bmod n$$

rank $g(\ell)$. ns

X:	0 1 2	3 4 5	6 7 8	9 10 11
g:	0 5 0	7 8 2	4 10	1 8 6

INOT: ~~no place to start~~ We need to start somewhere!

- ① Start with $h_1(\ell) = 1$ and $h_2(\ell) = 6$
Second ~~rank~~ $n(\ell) \rightarrow 1$ [Rank 1]

$$h(\ell) = g(h_1(\ell)) + g(h_2(\ell)) \bmod n$$

~~$$2 = g(1) + g(6) \bmod 9$$~~

$$h(t) = \phi(h_1(t)) + \phi(h_2(t)) \pmod{n}$$

Initially select node 0, set $\phi(0)$ need to start same.

~~$$\phi(t) = \phi(0) + \phi(t) \pmod{9}$$~~

~~$$\phi(N) = 0 + \phi(N) \pmod{9}$$~~

~~$$h(N) = \phi(h_1(N)) + \phi(h_2(N))$$~~

~~$$\text{We can set } \phi(1) = h(0) - \phi(0) \pmod{n}$$~~

~~$$h(t) = \phi(h_1(t)) + \phi(h_2(t)) \pmod{n}$$~~

Initially: select 0, and set $\boxed{\phi(0) = 0}$

~~$$h(0, 1) = \phi(\cancel{h_1(0)}) + \phi(1) \pmod{9}$$~~

~~$$S = 0 + \phi(1) \pmod{9}$$~~

~~$$\Rightarrow \phi(1) = S \pmod{9}$$~~

~~$$\Rightarrow \phi(1) = 5$$~~

Now \Rightarrow Continue from node(1), as you have info. motion about $\phi(1)$ there.

$$h(1, 0) = \phi(1) + \phi(0) \pmod{9}$$

$$0 = 5 + \phi(0) \pmod{9}$$

$$\Rightarrow \phi(0) = -5 \pmod{9}$$

$$\Rightarrow \phi(6) = 4$$

$$h(6, 4) = \phi(6) + \phi(4) \bmod 9$$

$$3 = 4 + \phi(4) \bmod 9$$

$$\phi(4) = -1 \bmod 9$$

$$\Rightarrow \phi(4) = 8$$

Let's continue with:

$$h(10, 1) = \phi(1) + \phi(10) \bmod 9$$

$$4 = *5 + \phi(10) \bmod 9$$

$$\phi(10) = * - 1 \bmod 9$$

$$\phi(10) = 8$$

(NO CONNECTED NODES LEFT)

\Rightarrow Selected a new node, listed (2) and set.

$$\phi(2) = 0$$

$$h(2, 7) = \phi(2) + \phi(7) \bmod 9$$

$$1 = *0 + \phi(7) \bmod 9$$

$$\Rightarrow \phi(7) = 1 \bmod 9$$

$$\phi(7) = 1$$

$$h(*7, 5) = \phi(-2) + \phi(5) \bmod 9$$

$$2 = 1 + \phi(4) \bmod 9$$

$$\phi(4) = *1$$

$$h(5, 3) = \phi(1) + \phi(3) \bmod 9$$

$$8 = *3 + \phi(3) \bmod 9$$

~~$$g(3) = 7 \pmod{9}$$

$$g(3) = 7$$~~

$$\boxed{g(3) = 7 \pmod{9}}$$

$$\Rightarrow \boxed{g(3) = 7}$$

< No mōles & leg >

Set $g(d) = 0$

$$h(8, 11) = g(8) + g(11) \pmod{9}$$

~~$$G = g(8) + g(11) \pmod{9}$$~~

$$g(11) = G \pmod{9}$$

$$\boxed{g(11) = G}$$

~~Set d~~

$$h(11, 9) = g(11) + g(9) \pmod{9}$$

$$7 = 6 + g(9) \pmod{9}$$

$$g(9) = 7 - 6 \pmod{9}$$

$$\boxed{g(9) = 1}$$

\Rightarrow NB. After finding mōles for all $g(t)$, we can similarly find $h(t)$ as:

$$\text{circled: } h(t) = g(h_1(t)) + g(h_2(t)) \pmod{9}$$

$$h(t_1, t_2)$$

Final value
connecting $h(t_1)$ and $h(t_2)$

THEOREM — Hash func to SAME BUCKET

For any two distinct keys x and y ,
the P. Func x hashes to the same bucket
giving $\Theta\left(\frac{1}{m}\right)$, where m is the
HASH TABLE's size.

CUCKOO HASHING \Rightarrow DYNAMIC HASH TABLE!

Innovative idea for solving HASH TABLES' creation, maintenance (i.e.: updated and re-order).

\rightarrow What? It combines pre-existing ideas into Ritterwurz to hash:

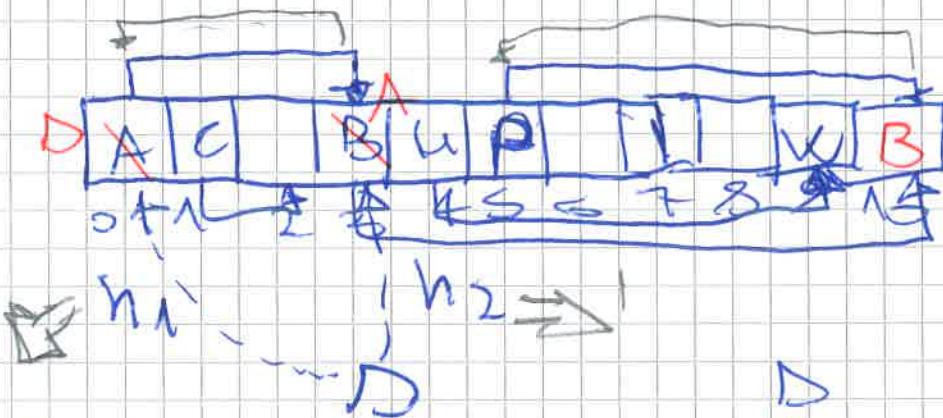
- DELENA - $O(n)$ SEARCH TIME (WORST CASE)
- "
- $O(n)$ UPDATE TIME (AMORTIZED COST)
- $O(n)$ SPACE

It is based on 2 hash FUNCTIONS: h_1 and h_2 that operate by the "PAIRS & 2" approach: when inserting a key, yet with a "TWIST".

\rightarrow Consider a HASH TABLE:

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 | 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 | 101 | 102 | 103 | 104 | 105 | 106 | 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 | 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 | 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 | 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 | 147 | 148 | 149 | 150 | 151 | 152 | 153 | 154 | 155 | 156 | 157 | 158 | 159 | 160 | 161 | 162 | 163 | 164 | 165 | 166 | 167 | 168 | 169 | 170 | 171 | 172 | 173 | 174 | 175 | 176 | 177 | 178 | 179 | 180 | 181 | 182 | 183 | 184 | 185 | 186 | 187 | 188 | 189 | 190 | 191 | 192 | 193 | 194 | 195 | 196 | 197 | 198 | 199 | 200 | 201 | 202 | 203 | 204 | 205 | 206 | 207 | 208 | 209 | 210 | 211 | 212 | 213 | 214 | 215 | 216 | 217 | 218 | 219 | 220 | 221 | 222 | 223 | 224 | 225 | 226 | 227 | 228 | 229 | 230 | 231 | 232 | 233 | 234 | 235 | 236 | 237 | 238 | 239 | 240 | 241 | 242 | 243 | 244 | 245 | 246 | 247 | 248 | 249 | 250 | 251 | 252 | 253 | 254 | 255 | 256 | 257 | 258 | 259 | 260 | 261 | 262 | 263 | 264 | 265 | 266 | 267 | 268 | 269 | 270 | 271 | 272 | 273 | 274 | 275 | 276 | 277 | 278 | 279 | 280 | 281 | 282 | 283 | 284 | 285 | 286 | 287 | 288 | 289 | 290 | 291 | 292 | 293 | 294 | 295 | 296 | 297 | 298 | 299 | 300 | 301 | 302 | 303 | 304 | 305 | 306 | 307 | 308 | 309 | 310 | 311 | 312 | 313 | 314 | 315 | 316 | 317 | 318 | 319 | 320 | 321 | 322 | 323 | 324 | 325 | 326 | 327 | 328 | 329 | 330 | 331 | 332 | 333 | 334 | 335 | 336 | 337 | 338 | 339 | 340 | 341 | 342 | 343 | 344 | 345 | 346 | 347 | 348 | 349 | 350 | 351 | 352 | 353 | 354 | 355 | 356 | 357 | 358 | 359 | 360 | 361 | 362 | 363 | 364 | 365 | 366 | 367 | 368 | 369 | 370 | 371 | 372 | 373 | 374 | 375 | 376 | 377 | 378 | 379 | 380 | 381 | 382 | 383 | 384 | 385 | 386 | 387 | 388 | 389 | 390 | 391 | 392 | 393 | 394 | 395 | 396 | 397 | 398 | 399 | 400 | 401 | 402 | 403 | 404 | 405 | 406 | 407 | 408 | 409 | 410 | 411 | 412 | 413 | 414 | 415 | 416 | 417 | 418 | 419 | 420 | 421 | 422 | 423 | 424 | 425 | 426 | 427 | 428 | 429 | 430 | 431 | 432 | 433 | 434 | 435 | 436 | 437 | 438 | 439 | 440 | 441 | 442 | 443 | 444 | 445 | 446 | 447 | 448 | 449 | 450 | 451 | 452 | 453 | 454 | 455 | 456 | 457 | 458 | 459 | 460 | 461 | 462 | 463 | 464 | 465 | 466 | 467 | 468 | 469 | 470 | 471 | 472 | 473 | 474 | 475 | 476 | 477 | 478 | 479 | 480 | 481 | 482 | 483 | 484 | 485 | 486 | 487 | 488 | 489 | 490 | 491 | 492 | 493 | 494 | 495 | 496 | 497 | 498 | 499 | 500 | 501 | 502 | 503 | 504 | 505 | 506 | 507 | 508 | 509 | 510 | 511 | 512 | 513 | 514 | 515 | 516 | 517 | 518 | 519 | 520 | 521 | 522 | 523 | 524 | 525 | 526 | 527 | 528 | 529 | 530 | 531 | 532 | 533 | 534 | 535 | 536 | 537 | 538 | 539 | 540 | 541 | 542 | 543 | 544 | 545 | 546 | 547 | 548 | 549 | 550 | 551 | 552 | 553 | 554 | 555 | 556 | 557 | 558 | 559 | 560 | 561 | 562 | 563 | 564 | 565 | 566 | 567 | 568 | 569 | 570 | 571 | 572 | 573 | 574 | 575 | 576 | 577 | 578 | 579 | 580 | 581 | 582 | 583 | 584 | 585 | 586 | 587 | 588 | 589 | 590 | 591 | 592 | 593 | 594 | 595 | 596 | 597 | 598 | 599 | 600 | 601 | 602 | 603 | 604 | 605 | 606 | 607 | 608 | 609 | 610 | 611 | 612 | 613 | 614 | 615 | 616 | 617 | 618 | 619 | 620 | 621 | 622 | 623 | 624 | 625 | 626 | 627 | 628 | 629 | 630 | 631 | 632 | 633 | 634 | 635 | 636 | 637 | 638 | 639 | 640 | 641 | 642 | 643 | 644 | 645 | 646 | 647 | 648 | 649 | 650 | 651 | 652 | 653 | 654 | 655 | 656 | 657 | 658 | 659 | 660 | 661 | 662 | 663 | 664 | 665 | 666 | 667 | 668 | 669 | 670 | 671 | 672 | 673 | 674 | 675 | 676 | 677 | 678 | 679 | 680 | 681 | 682 | 683 | 684 | 685 | 686 | 687 | 688 | 689 | 690 | 691 | 692 | 693 | 694 | 695 | 696 | 697 | 698 | 699 | 700 | 701 | 702 | 703 | 704 | 705 | 706 | 707 | 708 | 709 | 710 | 711 | 712 | 713 | 714 | 715 | 716 | 717 | 718 | 719 | 720 | 721 | 722 | 723 | 724 | 725 | 726 | 727 | 728 | 729 | 730 | 731 | 732 | 733 | 734 | 735 | 736 | 737 | 738 | 739 | 740 | 741 | 742 | 743 | 744 | 745 | 746 | 747 | 748 | 749 | 750 | 751 | 752 | 753 | 754 | 755 | 756 | 757 | 758 | 759 | 760 | 761 | 762 | 763 | 764 | 765 | 766 | 767 | 768 | 769 | 770 | 771 | 772 | 773 | 774 | 775 | 776 | 777 | 778 | 779 | 780 | 781 | 782 | 783 | 784 | 785 | 786 | 787 | 788 | 789 | 790 | 791 | 792 | 793 | 794 | 795 | 796 | 797 | 798 | 799 | 800 | 801 | 802 | 803 | 804 | 805 | 806 | 807 | 808 | 809 | 810 | 811 | 812 | 813 | 814 | 815 | 816 | 817 | 818 | 819 | 820 | 821 | 822 | 823 | 824 | 825 | 826 | 827 | 828 | 829 | 830 | 831 | 832 | 833 | 834 | 835 | 836 | 837 | 838 | 839 | 840 | 841 | 842 | 843 | 844 | 845 | 846 | 847 | 848 | 849 | 850 | 851 | 852 | 853 | 854 | 855 | 856 | 857 | 858 | 859 | 860 | 861 | 862 | 863 | 864 | 865 | 866 | 867 | 868 | 869 | 870 | 871 | 872 | 873 | 874 | 875 | 876 | 877 | 878 | 879 | 880 | 881 | 882 | 883 | 884 | 885 | 886 | 887 | 888 | 889 | 890 | 891 | 892 | 893 | 894 | 895 | 896 | 897 | 898 | 899 | 900 | 901 | 902 | 903 | 904 | 905 | 906 | 907 | 908 | 909 | 910 | 911 | 912 | 913 | 914 | 915 | 916 | 917 | 918 | 919 | 920 | 921 | 922 | 923 | 924 | 925 | 926 | 927 | 928 | 929 | 930 | 931 | 932 | 933 | 934 | 935 | 936 | 937 | 938 | 939 | 940 | 941 | 942 | 943 | 944 | 945 | 946 | 947 | 948 | 949 | 950 | 951 | 952 | 953 | 954 | 955 | 956 | 957 | 958 | 959 | 960 | 961 | 962 | 963 | 964 | 965 | 966 | 967 | 968 | 969 | 970 | 971 | 972 | 973 | 974 | 975 | 976 | 977 | 978 | 979 | 980 | 981 | 982 | 983 | 984 | 985 | 986 | 987 | 988 | 989 | 990 | 991 | 992 | 993 | 994 | 995 | 996 | 997 | 998 | 999 | 1000 | 1001 | 1002 | 1003 | 1004 | 1005 | 1006 | 1007 | 1008 | 1009 | 1010 | 1011 | 1012 | 1013 | 1014 | 1015 | 1016 | 1017 | 1018 | 1019 | 1020 | 1021 | 1022 | 1023 | 1024 | 1025 | 1026 | 1027 | 1028 | 1029 | 1030 | 1031 | 1032 | 1033 | 1034 | 1035 | 1036 | 1037 | 1038 | 1039 | 1040 | 1041 | 1042 | 1043 | 1044 | 1045 | 1046 | 1047 | 1048 | 1049 | 1050 | 1051 | 1052 | 1053 | 1054 | 1055 | 1056 | 1057 | 1058 | 1059 | 1060 | 1061 | 1062 | 1063 | 1064 | 1065 | 1066 | 1067 | 1068 | 1069 | 1070 | 1071 | 1072 | 1073 | 1074 | 1075 | 1076 | 1077 | 1078 | 1079 | 1080 | 1081 | 1082 | 1083 | 1084 | 1085 | 1086 | 1087 | 1088 | 1089 | 1090 | 1091 | 1092 | 1093 | 1094 | 1095 | 1096 | 1097 | 1098 | 1099 | 1100 | 1101 | 1102 | 1103 | 1104 | 1105 | 1106 | 1107 | 1108 | 1109 | 1110 | 1111 | 1112 | 1113 | 1114 | 1115 | 1116 | 1117 | 1118 | 1119 | 1120 | 1121 | 1122 | 1123 | 1124 | 1125 | 1126 | 1127 | 1128 | 1129 | 1130 | 1131 | 1132 | 1133 | 1134 | 1135 | 1136 | 1137 | 1138 | 1139 | 1140 | 1141 | 1142 | 1143 | 1144 | 1145 | 1146 | 1147 | 1148 | 1149 | 1150 | 1151 | 1152 | 1153 | 1154 | 1155 | 1156 | 1157 | 1158 | 1159 | 1160 | 1161 | 1162 | 1163 | 1164 | 1165 | 1166 | 1167 | 1168 | 1169 | 1170 | 1171 | 1172 | 1173 | 1174 | 1175 | 1176 | 1177 | 1178 | 1179 | 1180 | 1181 | 1182 | 1183 | 1184 | 1185 | 1186 | 1187 | 1188 | 1189 | 1190 | 1191 | 1192 | 1193 | 1194 | 1195 | 1196 | 1197 | 1198 | 1199 | 1200 | 1201 | 1202 | 1203 | 1204 | 1205 | 1206 | 1207 | 1208 | 1209 | 1210 | 1211 | 1212 | 1213 | 1214 | 1215 | 1216 | 1217 | 1218 | 1219 | 1220 | 1221 | 1222 | 1223 | 1224 | 1225 | 1226 | 1227 | 1228 | 1229 | 1230 | 1231 | 1232 | 1233 | 1234 | 1235 | 1236 | 1237 | 1238 | 1239 | 1240 | 1241 | 1242 | 1243 | 1244 | 1245 | 1246 | 1247 | 1248 | 1249 | 1250 | 1251 | 1252 | 1253 | 1254 | 1255 | 1256 | 1257 | 1258 | 1259 | 1260 | 1261 | 1262 | 1263 | 1264 | 1265 | 1266 | 1267 | 1268 | 1269 | 1270 | 1271 | 1272 | 1273 | 1274 | 1275 | 1276 | 1277 | 1278 | 1279 | 1280 | 1281 | 1282 | 1283 | 1284 | 1285 | 1286 | 1287 | 1288 | 1289 | 1290 | 1291 | 1292 | 1293 | 1294 | 1295 | 1296 | 1297 | 1298 | 1299 | 1300 | 1301 | 1302 | 1303 | 1304 | 1305 | 1306 | 1307 | 1308 | 1309 | 1310 | 1311 | 1312 | 1313 | 1314 | 1315 | 1316 | 1317 | 1318 | 1319 |<th
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Insertion: Ex: We want to insert $\langle 5 \rangle$ at position \rightarrow Operate by the "Power of 2" ARGUMENT



~~If both sets check $h_1(k)$ and $h_2(k)$. If either set is empty insert D there.~~

~~If both sets are filled with another value, then \Rightarrow find element in $h_1(k)$~~

~~① Insert D into set 1 & set 2. Kick out A from 1~~

~~→ Kick A out of set ① \Rightarrow Insert D~~

~~into set ① \Rightarrow ~~Freeze dog~~ from ① to ② \Rightarrow Kick B from~~

~~set ② \Rightarrow Insert A into set ②~~

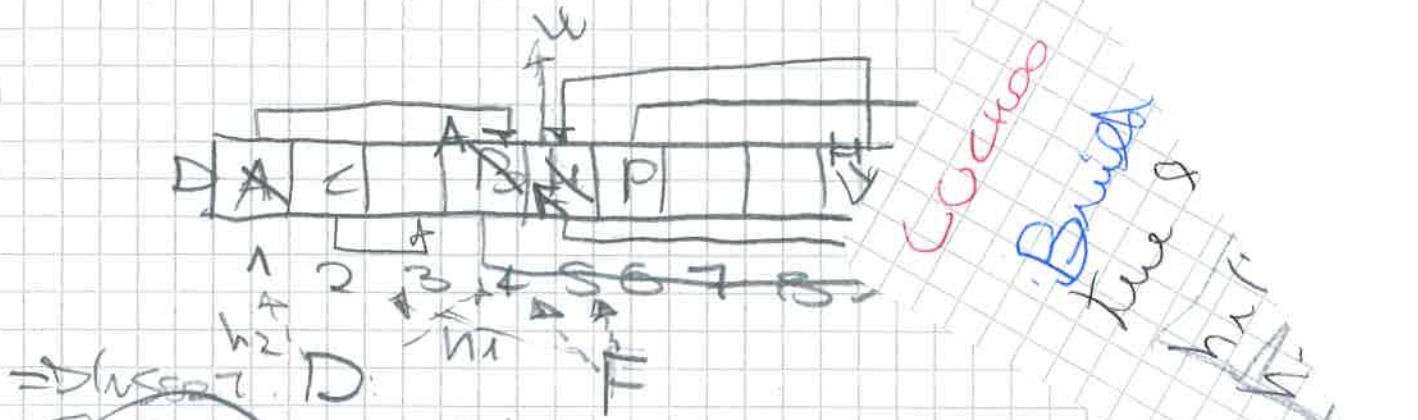
~~Very movements \Rightarrow ~~Freeze dog~~ from ② to ③
 \Rightarrow Insert B into ③, as it's empty.~~

[We keep following dogs till we find an empty slot] \Rightarrow All cells are full,

then we can't ~~insert~~

~~insert anything~~

Book
EXAMPLE - Cuckoo Hashing's "



$\Rightarrow \text{INSERT } D$: $h_1(D) = 4, h_2(D) = 1 \Rightarrow D$ Both are free.

\Rightarrow The pick ① for eviction.
A

$\Rightarrow \text{INSERT } F$:

$$h_1(F) = 4; h_2(F) = 5$$



$$h_1(F) = 4; h_2(A) = 5$$

Covers "LOOLED" Evictions
Knot stop upon Reaching ⑤
Again. At that point, we
also check $h_1(K)$ for insertion

\Rightarrow After ~~inserting~~ inserting elements, we need
to ~~remove~~ ~~remove~~ INSERT the slots

/EVICION: Firstly call h_1 . After
fully processing it, we will have one
extra element | also shade h_2 .

Cuckoo Hashing - CLASS EXERCISE

Build a Cuckoo Hash Table based on the following two hash functions.

$$h_1(k) = k \bmod 17 \quad = D(m=17)$$

$$h_2(k) = 2k \bmod 17$$

1	18	3	21	5	20		22	23							
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

- $h_1(k)$ taken, h_2 not taken \Rightarrow Insert in $h_2(k)$.
- $h_1(k), h_2(k) \Rightarrow$ Need to check w/ $h_1(k)$.

$$\begin{array}{r} 76 \\ 34 \\ \hline 12 \end{array}$$

$$S = \{ \cancel{1}, \cancel{5}, \cancel{18}, \cancel{3}, \cancel{22}, \cancel{20}, \cancel{21}, \cancel{23}, \cancel{19} \}$$

$$\begin{array}{r} 42 \\ 30 \\ \hline 12 \end{array}$$

$$h_1(19) = 2$$

$$h_2(19) = 4$$

$$\begin{array}{|c|} \hline h_1(k) \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline h_2(k) \\ \hline \end{array}$$

$$o(h_1(k), h_2(k))$$

18	18	19	*												
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

After

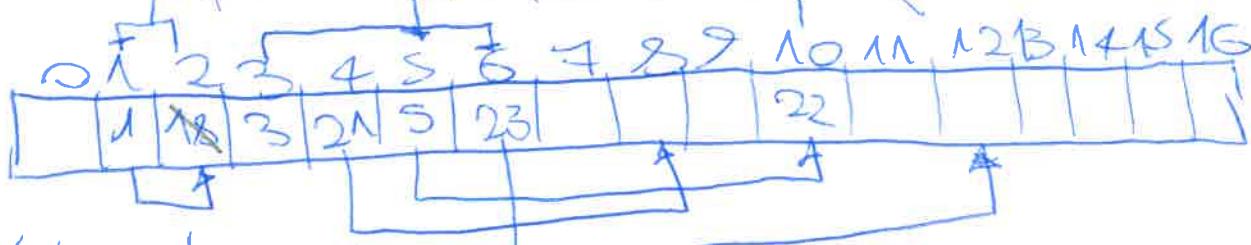
swap

NB: When performing "cyclic" evictions, we need to check h_2

\Rightarrow Also, we need to update single direction!

CWEKOO WASHING - EXERCISE WITH FRANCESCO

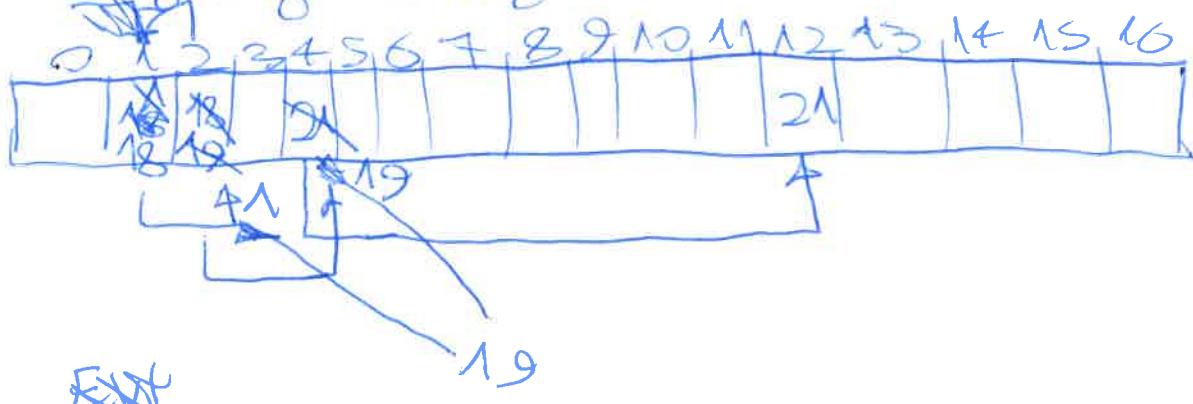
$$S = \{x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{12}, x_{13}, x_{14}, x_{15}, x_{16}\}$$



$$h_1(k) = k \bmod 17$$

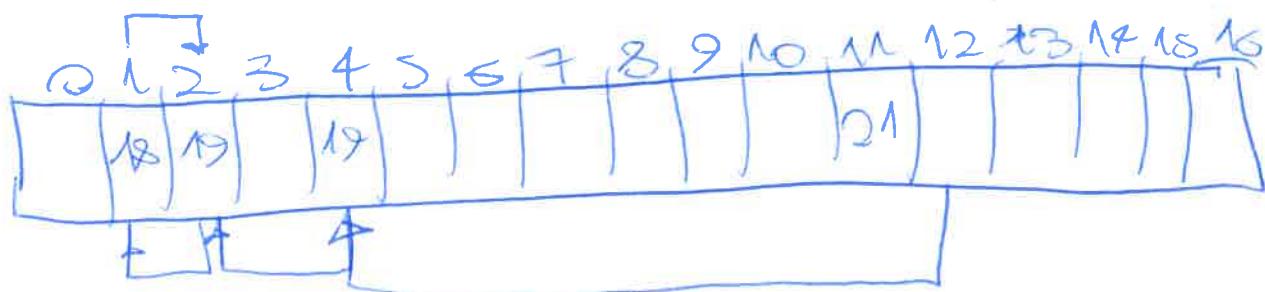
$$h_2(k) = 2k \bmod 17$$

Everything straight till 19 \Rightarrow Now Error



~~Excess~~

And now INVERT the EDGES:



RESOLVING FROM TABLE.

18	1	3	8	19	5	20	21	22	23	
0	1	2	3	4	5	6	7	8	9	10
1	1	4	4				1	4	4	4

THEOREM - Hashing to the SAME BUCKET

INTUITION: Hashing:

For any two DISTINCT keys x and y , the $P\{x \text{ hashes to the same bucket as } y\} = O(\frac{1}{m})$

Proof:

If x and y are in THE SAME Bucket, Then
There is a PATH of some length L between one
node in $\{h_1(x), h_2(x)\}$ and one node in $\{h_1(y), h_2(y)\}$

To be shown by PROV. THEOREM that
All entries in $\{y\}$ only contain $L > 1$

If $m \geq 2 \cdot C \cdot n$, then:

$$P\{i \neq j\} = \frac{C^L}{m}$$

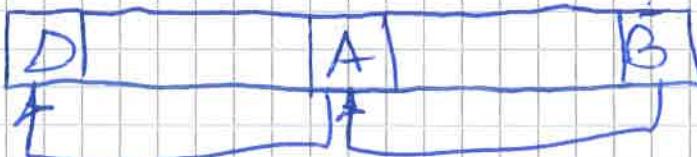
$$\Rightarrow 4 \cdot \sum_{L=1}^{+\infty} C^L = \frac{4}{m} \sum_{L=1}^{+\infty} \left(\frac{1}{2}\right)^L = \frac{4}{m} \cdot \frac{\frac{1}{2}}{1-\frac{1}{2}}$$

$$= \frac{4}{m} \cdot \frac{\frac{1}{2}}{C-1} = \frac{4}{C \cdot m} = O\left(\frac{1}{m}\right)$$

Q.E.D.

N.B.: After inserting new values, we need to INVEST the EXISTS' DIRECTION of the ones percolated back following an EVACUATION.

Ex.:



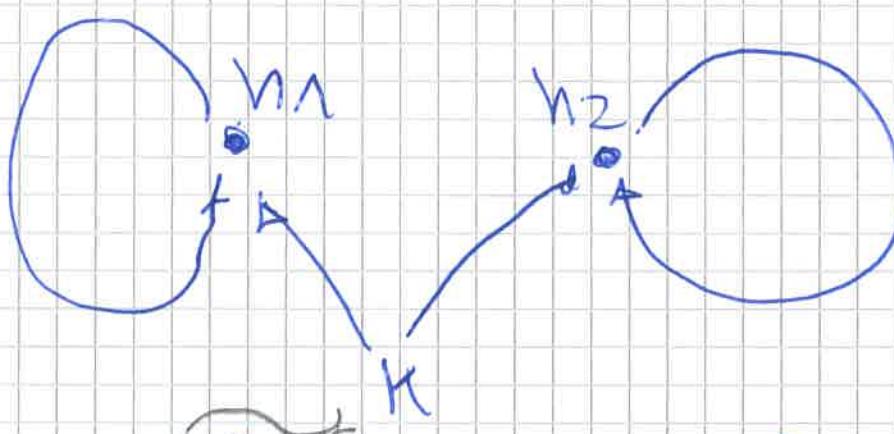
CYCLES? How to deal with them:

↳ It can happen that both functions lead to the creation of cycles among keys!

- ↳ • 1 CYCLE: N.P! \rightarrow Periodic pointers.
- 2+ CYCLES: issues! \rightarrow By REHASHING small elements

Need to pick new hash functions and
rehash all the dictionary keys
(however, keys occurs with bounded P.)

↳ lot of rehashing can be amortized by cost of O/N for inserting.



PERSOF. Having one cycle is very improbable
(i.e. The O(N) cost of rehashing is
AMORTIZED \rightarrow by insertion taking O(1)
Time).

Ex. The 20,000 € of losing a new
VNL Pdo to amortize over 10 years.

THEOREM: $L = \text{Path length between } i \text{ and } j.$

If $m = \text{array size} = \# \text{NODES}$

$n = \# \text{edges} = \# \text{EDGES}$

THE LENGTH of a PATH:

Any entries in Δ and any constant $C > 1$.
↳ Markable

If $m \geq 2 < n$, (at most S_m edges)

$$P\{\Delta_{ij} \leq L\} \leq C^{-L} = \frac{1}{m \cdot C^L}$$

i.e.,

The Probability To peripheral a path gets always smaller and smaller (exponentially).

PROOF - By induction over (i) and (j) :

• BASE CASE: $L = 1$ (i.e. one edge between i and j)

Path length 1 (1 edge)
 $\frac{1}{m^2}$ (number of edges among m^2 vertices)

$$P\{\Delta_{ij} = 1\} = \sum_{k=1}^{m^2} \frac{2}{m^2} = \frac{2m}{m^2} < \frac{m}{2 \cdot m}$$

edge between (i) and (j)

$$\text{hyp: } m \geq 2 \cdot \ln \frac{1}{C \cdot m} \Rightarrow n \leq \frac{m}{22}$$

$$P\{\Delta_{ij} \leq 1\} \leq \frac{1}{cm}$$

• INDUCTIVE CASE: $L > 1$

Assume the bound holds for $(L-1)$. We want to show that it holds for L as well.

~~SITUATION VISUALIZED~~



Known by

INDUCTIVE hypothesis

, namely!

RES

$$\sum_{L=1}^{L-1} \frac{1}{m} = \frac{L-1}{m}$$

$$\cdot \frac{1}{m}$$

INDUCTIVE

HYP

$$(i \rightarrow j)$$

remaining

EDGE

$$m \rightarrow j$$

C.Q.F.D.

~~$\sum_m \sum_{L=1}^{L-1} \frac{1}{m}$~~

$$\# = \frac{L-1}{m}$$

(NB: By INDUCTIVE hypothesis, we had that

$$P\{i \rightarrow j\} \leq \frac{1}{m}$$

PROBABILITY OF A CYCLE EXISTENCE:

We now want to prove that:

$$P\{f^{\circ L} = i\}$$

$$\# \sum_{L=1}^{L-1} P\{i \rightarrow i\} = \sum_{L=1}^{L-1} \frac{1}{m \cdot c^L}$$

(There exists)
a cycle in

Consider all the possible
paths length L having
to a cycle

$$= \frac{1}{m} \sum_{L=1}^{L-1} \left(\frac{1}{c}\right)^L$$

$$= \frac{1}{m} \cdot \cancel{\frac{(1)}{(1-1/c)}} \frac{1}{m^{L-1}} = \frac{1}{m^{L-1}}$$

$$= P\{\text{cycle in } i\} \leq \frac{1}{m} \cdot \frac{1}{c-1}$$

$$P\{\text{CYCLE EXISTENCE}\} \leq m \cdot \frac{1}{m} \cdot \frac{1}{c-1}$$

$$P\{\text{CYCLE}\} \leq \frac{1}{c-1}$$

Upper bound to
the P. of
unsuccessful
invention,

Corollary 1:

We have $\Omega(n)$ first.

$$P\{ \text{cycle} \} \leq \frac{1}{C-1}$$

$$\Rightarrow \text{We would like } \frac{1}{C-1} = \frac{1}{2} \Rightarrow 2 = C-1 \Rightarrow C = 3$$

Set $C \geq 3$, and we have that:

$$m \geq 2C \cdot n$$

$$\Rightarrow m \geq 6 \cdot (1+\varepsilon) n$$

Taking a **subset TABLE** of this size, the P. to have a cycle in the subset graph of the original dictionary is $\leq \frac{1}{2}$

$$P\{ \text{cycle} \} \leq \frac{1}{2}$$

\Rightarrow A constant #REHASHES is enough to ensure that the insertion of Σn keys in a dictionary of size $n \Rightarrow$ One rehashing of tables $O(n) \Rightarrow O(1/\varepsilon)$ per insertion (amortized).

Corollary 2:

By setting $C \geq 2$ and taking a subset table of size $m \geq 2cn(1+\varepsilon)$, the cost for inserting Σn keys in a DICTIONARY of size n is constant AMORTIZED (AMORTIZED over the cost of $\Theta(n)$ insertions a DE-HASH).

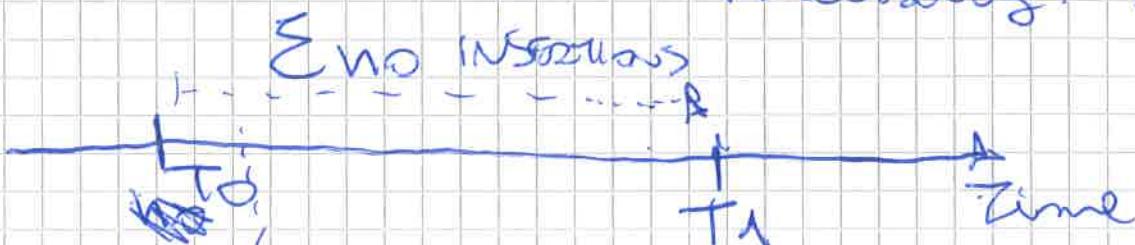
~~Settler~~

We want to show that, if the ~~cost~~

Keys $\leq (1+\epsilon)N \Rightarrow$ Theorem holds!

Let's consider a hash TABLE for ~~the BTB~~

$m = 6 \cdot \epsilon \cdot k_n$ (i.e. larger than necessary!)



Let's consider:

$P\{ \text{Sail at } T_0 \} < P\{ \text{Sail at } T_0 + 1 \}$

\Rightarrow {table reconstruction} \leftarrow {table reconstruction}
at T_0 at $T_0 + 1$

→ Always starting P. to sail

\Rightarrow Let's now consider the time interval
 T_1 .

$$\{P\{ \text{Sail at } T_1 \}\} \leq \frac{1}{2}$$

[On average two reconstructions are enough at T_1]

\Rightarrow The ADVERTISED COST is constant!
Because we have

~~Since two reconstructions~~

\Rightarrow Each insertion would take $O(\frac{1}{2})$

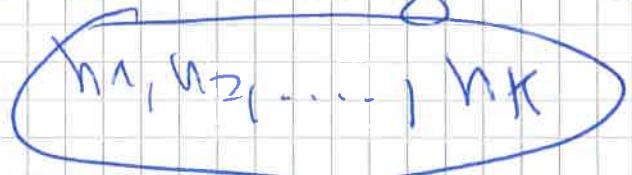
~~GLOBAL REBUILDING TECHNIQUE~~

Analogously to insertion with chaining, if the size of the dictionary becomes too small compared to the size of the hash table, a new smaller hash table ~~is~~ is created;
[Same ~~space~~ space]

If the size of the ~~dictionary~~ ~~is too large~~ hash table is filled up with all values of the dictionary, a new larger hash table is re-created (doubled, or by a constant factor).

~~K HASH FUNCTIONS.~~ [MULTI-CAPTE]

Consider K hash functions instead of 2.



We have a tree-like approach for having cell (key) => COMPLEX VISUALIZATION

~~BUCKET-LIKE APPROACH.~~

Every cell has a ~~Bucket~~ (we no longer have cells!)



We can calculate a ~~LOAD FILLING FACTOR~~

BLOOM FILTERS

Data structure used to store a very large universe of keys, hence occupying a large space [Ex: URLs managed by crawler in search engines]. Very long!
→ Need to keep them in memory!

Idea: Do not store strings themselves, but a simple and randomized representation of them (Fingerprint) in a **BINARY VECTOR**.

→ ~~DEAUBACK~~: One - table - error, yet error size decreases exponentially with the size of the FINGERPRINT of keys.

↳ $K \in D \Rightarrow$ Answer is **yes**
[No Error]

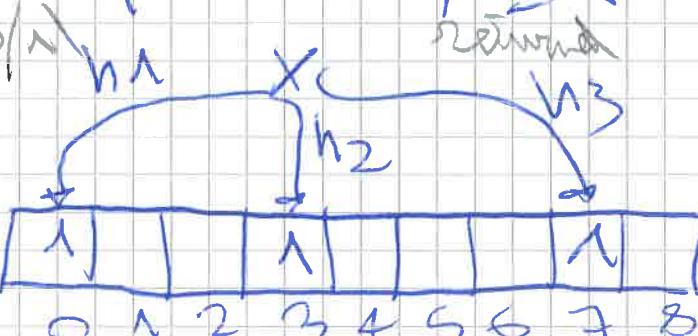
$K \notin D \Rightarrow$ Answer is ~~→ (YES/NO)~~
YES
↓ NO

→ There may occur **FALSE POSITIVES** in the search, in case elements are NOT in the dictionary D .

→ SIMPLE, RANDOMIZED APPROXIM. THE SEARCH onto HASH FUNCTIONS:



Efficiency:

- INSERT(x) \Rightarrow Set all lists produced by hash functions h_1, \dots, h_n to \square .
 $\mathcal{O}(n)$
- SEARCH(x) \Rightarrow YES \rightarrow All lists of $B[h_i(x)]$ are λ .
 $\mathcal{O}(n)$


My FALSE POSITIVES are POSSIBLE!

Search. Re-compute hash function by x and see if those lists are λ .
 FALSE: However, those lists may have been set to λ by another hash function!
 TRUE: (x) λ

\Rightarrow The P. of error depends on K
 (i.e., the hash functions used)

ESTIMATION of the P. of FALSE

~~FALSE~~

~~Ex: Insert x in P always yields a FALSE Positive~~

~~To calculate P as follows: Count the cases in which we have a FALSE Positive~~
 ~~$P = \{$ cells after cell x is \square after inserting n keys $\}$~~

~~= P independent hash functions $H_i(k)$ return an entry different from λ~~

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

ESTIMATION OF THE P. OF FALSE IN A BLOOM FILTER

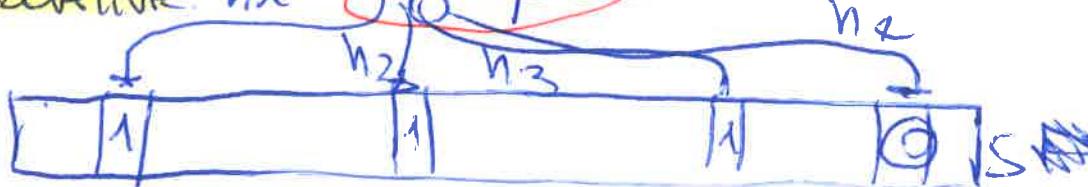
TRUE

NEGATIVE: h_1

~~h₁ & S~~

Correct! \Rightarrow h₁ is

not detected
to lie
in ΔS



~~h₁ & S~~

FALSE POSITIVE:



Wrong: h_1 is
detected to lie in
 S , though it
is actually
not there.

This bit may have been set to 1 because of other hashes of other keys.

$P\{$ FALSE-POSITIVE ERROR $\}$ = ?
in a Bloom FILTER

PROOF: ($\text{Estimate } P. \text{ of FALSE POSITIVE in a Bloom FILTER}$)

$P\{$ The insertion of a key ~~left~~ left a BIT-END $\}$
 $B[i] \Rightarrow \boxed{0}_i \} = P\{ K \text{ independent hash functions returned an entry } \neq i \}$
because $\boxed{0}_i$
 $= \left(\frac{m-1}{m} \right)^K = \left(1 - \frac{1}{m} \right)^K$

After inserting n -many keys, the P.
that $B[i]$ is still 0 is:

$$= \left(1 - \frac{1}{m}\right)^{n \cdot K}$$



$P\{\text{Bit } B[i] = 1 \text{ after inserting } K \text{ items}\} = \left(1 - e^{-\frac{n \cdot K}{m}}\right)$

$$= \left(1 - \frac{1}{m}\right)^{\frac{n \cdot K \cdot m}{m}} \approx e^{-\frac{n \cdot K}{m}}$$

$\times \notin D_i$, answering to $\exists x \in m \rightarrow 1$

$P\{\text{FALSE-POSITIVE Error}\} = P\{\text{all } K \text{ bits checked}$

$= P\{\text{all } K \text{ bits checked for a key not in the current dictionary are set to } 1\}$

$$= P\{\forall x : B[n_i(x) = 1]\} = \left(1 - e^{-\frac{n \cdot K}{m}}\right)^K$$

$$\Rightarrow \text{Perr} = \left(1 - e^{-\frac{n \cdot K}{m}}\right)^K$$

Whether the error probability depends on:

- K = # hash functions (the higher, the lower)
- m = # bits in the ~~binary array~~ (the more, the better)
- n = # Dictionary Keys (the fewer, the better)

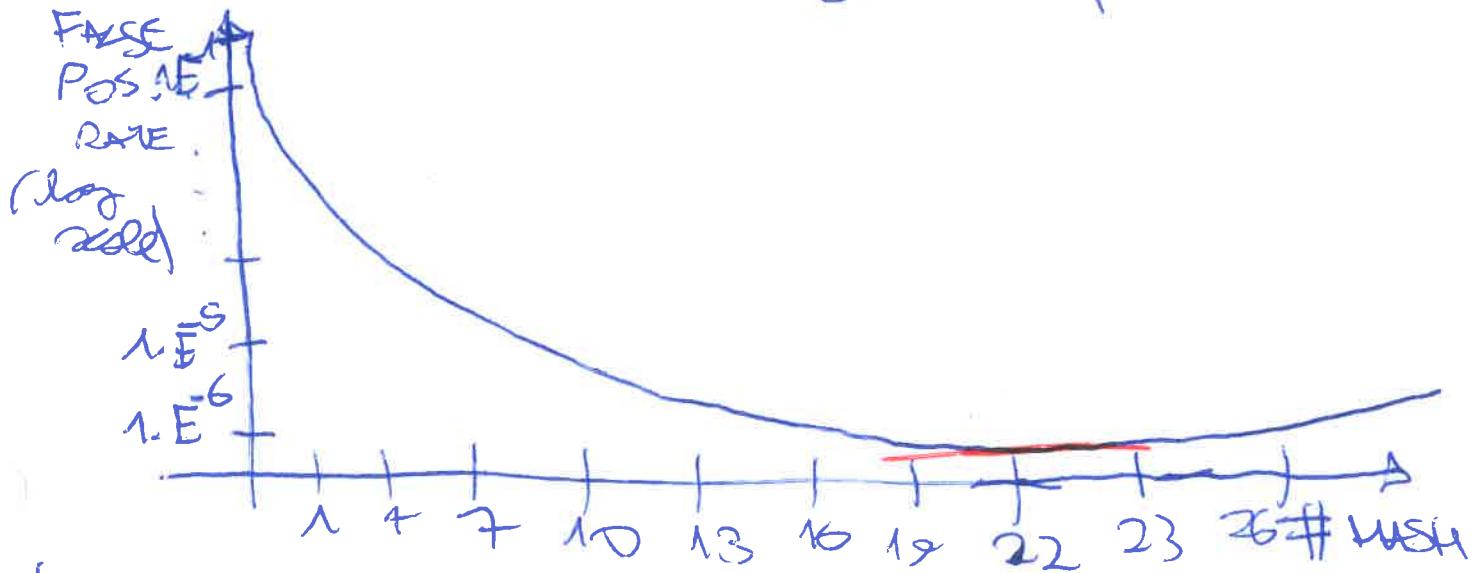
NB: $\beta = \frac{m}{n} = \text{Ava. \# Bits per Dictionary key allocated in } B / \text{fingerprint size}$

The larger β , the smaller the **Error Probability**, but the larger the **Space** allocated

~~OPTIMAL k~~

OPTIMAL VALUE FOR k ~~(k=2)~~

$k = \# \text{HASH FUNCTIONS}$ ~~to use~~ ~~needed~~



If we fix m and n ~~the # keys.~~, ~~functions~~
BLOOM FILTER'S SIZE.

what is the OPTIMAL value for k ?

(k^* : the one value minimizing PERR)

We know that:

$$P\{x \notin D, \text{answering "YES"}\} = (1 - e^{-\frac{nk}{m}})^k$$

OPTIMAL k: ~~# hash functions~~ We compute $\frac{d}{dk} P(k=0)$ to get

$$k = \frac{m}{n} \cdot \ln 2 = \left(\frac{1}{2}\right)^{\frac{m}{n} \cdot \ln 2} \text{ the MINIMUM.}$$

~~Large k \Rightarrow more 1s
Small k \Rightarrow more 0s~~ $\begin{matrix} \text{N.B.} \\ \text{B.} \end{matrix}$

\Rightarrow For OPTIMAL k \Rightarrow Same # 0s & 1s in B.

$$\text{PERR} = (0.6180)^{\frac{m}{n}} \text{ B!}$$

for k^* optimal

SPACE of Bloom Filter:

Lower-Bounded by:

$$m \geq \frac{n \cdot \log_2 \left(\frac{1}{\epsilon} \right)}{\ln 2} \approx 1.44 \cdot n \cdot \log_2 \left(\frac{1}{\epsilon} \right)$$

\Rightarrow SPACE = $O(n)$ = OPTIMAL!

(likely)

→ To store n keys.

UNIVERSAL HASH FUNCTIONS - CLASS

EXERCISE:

$$S = \{N, \$, \#, \%, \&, \&, \&, \&, \&\}$$

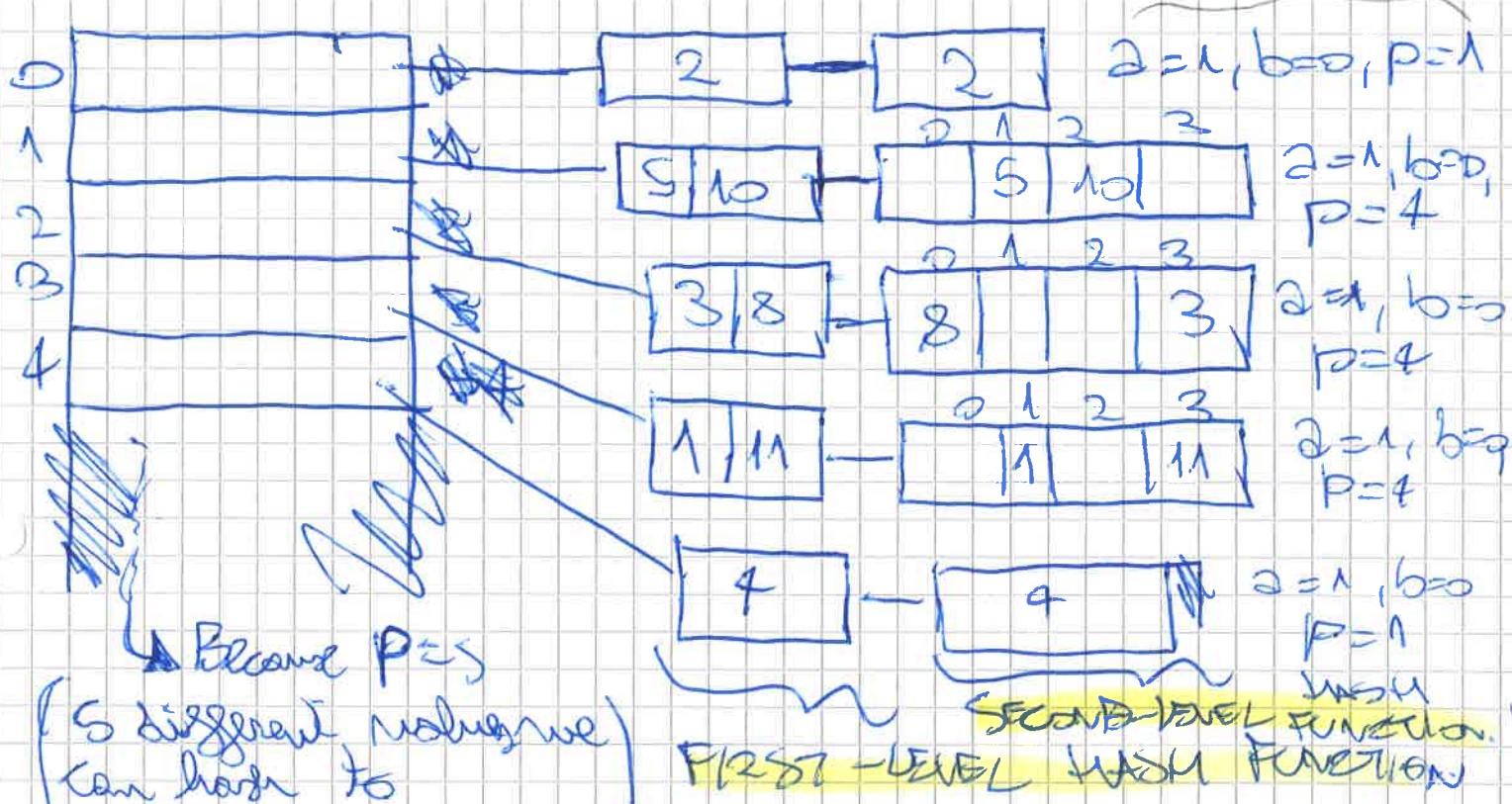
$$h_{a,b}(x) = ax + b \bmod p$$

First-level hash function:

~~$$h_1(x) = x \bmod 14$$~~

$$h_{1,0} = x \bmod 14$$

$$h_{2,1} = 2x + 1 \bmod 5$$



Is it perfect? (i.e.: no collisions).

$$\sum_{i=1}^5 (n_i)^2 = 14 < 2n - 16$$

n_i # Elements in bucket i is 16

[Otherwise, would need a different hash function!]

EX: Want to check if a number no contains in the hash table.

$$\rightarrow x = 7$$

$$1) n_1 = 2x + 1 \text{ mod } 5$$

$n_1 = 0 \rightarrow$ Go to Bucket 0.

$$2) n_2 = x \text{ mod } 1$$

NOT PRESENT!

In order to pick a UNIVERSAL HASH FUNCTION, generally 2 trials are sufficient.

$$\text{hash}(d = 2x + b \text{ mod } p)$$

where $p \geq m$

& PRIME NUMBER.

MINIMAL ORDERED PERFECT HASH FUNCTION.

$$S = \{aa, ab, ba, bc, ca, db\}$$

And we are given:

$$h_1(x, y) = \text{rank}(x) \cdot \text{rank}(y) \bmod m$$

$$h_2(x, y) = \text{rank}(x) + \text{rank}(y) \bmod m$$

~~$$h_1 = h_2$$~~



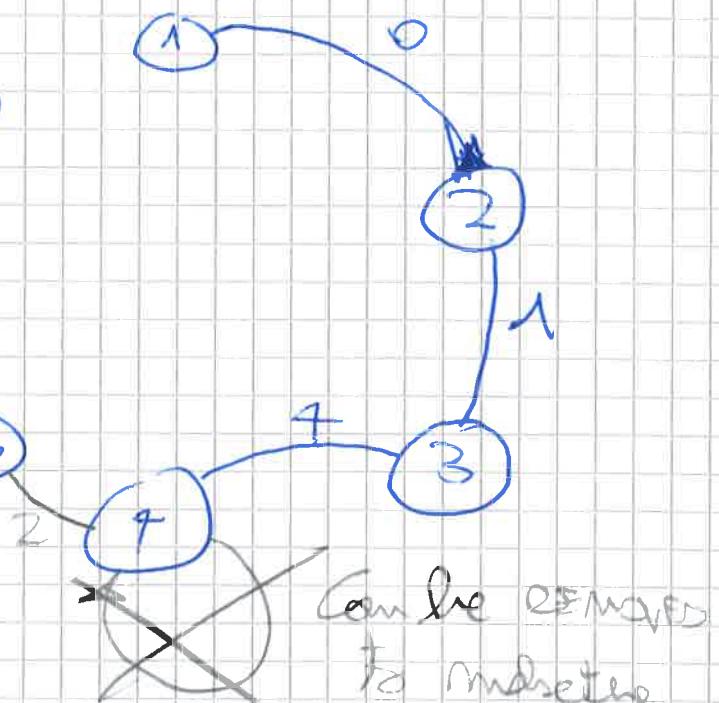
And also: $m = 11$,
 $n = 6$

We compute the respective hash functions

	h_1	h_2	rank
aa	2	1	0
ab	2	3	1
ba	4	5	2
bc	6	5	3
ca	3	4	4
db	8	6	5

If

We notice that ~~we have~~ we have as many nodes as the # different hash values.



→ We now need to find values for g_2 :

g	x
0	1
0	2
1	3
3	4
5	5
4	6
1	8

$$\bullet h(1,2) = g(1) + g(2) \bmod 6$$

$$\Rightarrow g = 0 + g(2) \bmod 6$$

$$g(2) = 0 \bmod 6$$

$$\bullet h(2,3) = g(2) + g(3) \bmod 6$$

$$1 = 0 + g(3) \bmod 6$$

$$g(3) = 1 \bmod 6$$

$$\bullet h(3,4) = g(3) + g(4) \bmod 6$$

~~\bullet~~ $4 = 1 + g(4) \bmod 6$

$$g(4) = 3 \bmod 6$$

$$\bullet h(4,5) = g(4) + g(5) \bmod 6$$

~~\bullet~~ ~~$5 = 3 + g(5) \bmod 6$~~

$$2 = 3 + g(5) \bmod 6$$

$$g(5) = -1 \bmod 6$$

$$g(5) = 5$$

$$\bullet h(5,6) = g(5) + g(6) \bmod 6$$

~~\bullet~~ $3 = 5 + g(6) \bmod 6$

$$g(6) = -2 \bmod 6$$

$$g(6) = 4 \bmod 6$$

$$\bullet h(6,7) = g(6) + g(7)$$

$$5 = 4 + g(7) \bmod 6$$

$$g(7) = 1$$

N.B.: This works only for strings actually inserted in the base table!

Not for other things.

BLOOM FILTER - CLASS EXERCISE:

1) Given $m = 2^{20}$ and $n = 2^{16}$
 75 keys to insert
 Bloom filter's size

→ Compute the OPTIMAL # hash functions to use & the error.

$$k_{opt} = \text{OPTIMAL } \# \text{hash functions} = \frac{m \ln 2}{n}$$

$$= \frac{2^{20}}{2^{16}} \ln 2 = 2^4 \ln 2$$

$$\text{Prob. Error} = e^{-k_{opt}^2} = (0.6185)^{\frac{m}{n}} = (0.6185)^{\frac{2^4}{2^{16}}} = (0.6185)^{\frac{1}{16}}$$

2) Compute size of Bloom Filter (m) to guarantee an error of $\epsilon = 2^{-20}$ by using an optimal # hash functions & knowing we have to insert $n = 2^{16}$ keys.

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

$$\epsilon = (0.6185)^{\frac{m}{n}} \Rightarrow 2^{-20} = (0.6185)^{\frac{m}{2^{16}}} \Rightarrow 2^{20} = (0.6185)^{\frac{m}{2^{16}}}$$

→ After finding (m), we can then find k_{opt} .

$$2^{20} = (0.6185)^m$$

$$m = 5.78 \quad k_{opt} = \frac{5.78}{2^{16}} \ln 2 = ?$$

$$\text{error} = \left[1 - e^{-\frac{Kn}{m}} \right]^K$$

SEARCHING STRINGS BY PREFIX

[CHAPTER 9]

PREFIX SEARCH

PROBLEM STATEMENT.

INPUT: • Dictionary D consisting of n strings.

• Query P = A string containing a set of characters.

OUTPUT: Strings of D that have P as a PREFIX.

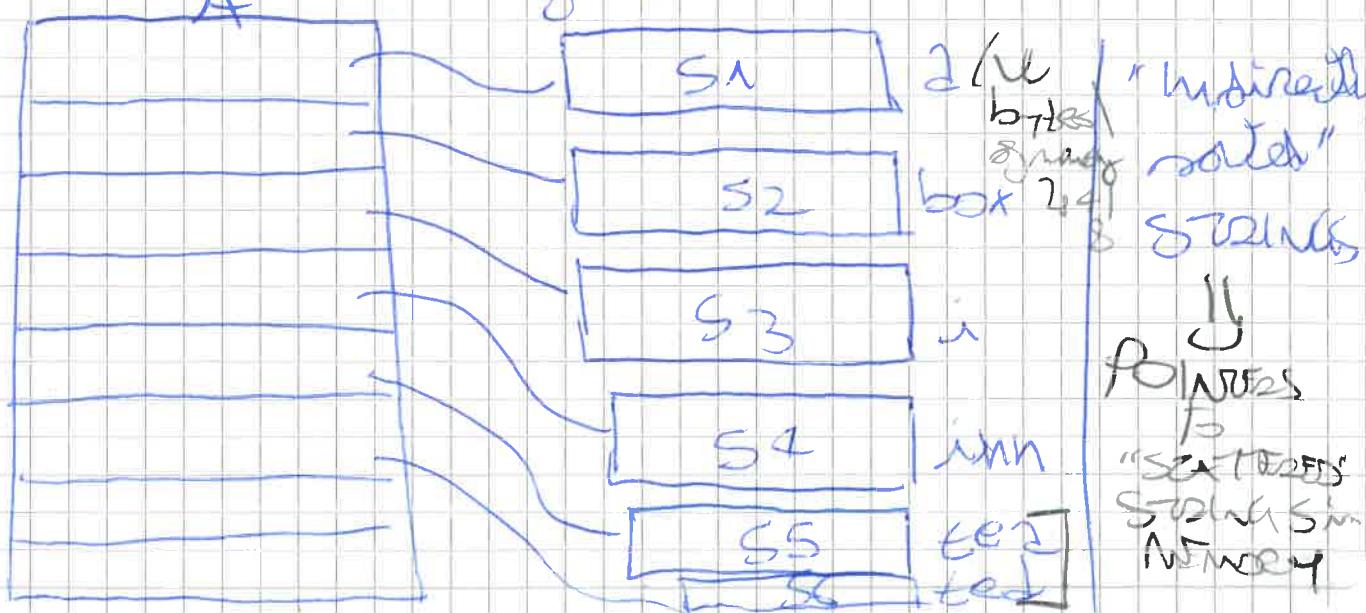
n = #STRINGS in D .

N = Total size of D .

$P = [1, P]$ = Array of $|P|$ characters.

① SOLUTION - ARRAY & STRING POINTERS.

We take an INPUT: Array A of pointers stored in arbitrary locations on disk.



A trie is the result of calling malloc for n strings.

↳ Arbitrary strings' allocation, instead memory management is delegated to the SYSTEM.

BASIC APPROACH'S IDEA:

Perform a simple binary search of P over the strings' PREFIXES in the dictionary.

COMPARISON STEPS = $O(\log_2 n)$

COST OF 1 STRING COMPARISON = $O(p)$ in worst case, from all the prefixes

∴ PREFIX BINARY SEARCH COST = $O(P \cdot \log_2 n)$

SPACE = $O(N + n)$

Ex. $P = \text{te}$

∴ We are interested in finding tea , ted .

Goal: Identifying Block's starting and ending position of the blocks of strings where tea and ted are located.
[i.e. of all strings having that prefix]

WE EXPLAIN 2 MAIN PROPERTIES for this
BASIC APPROACH:

- 1) All deletion array strings ~~prefixed by~~ by pattern P are contiguous if they are lexicographically sorted, so their pointers ~~occupy~~ occupy a continuous subarray, say $A[l, r]$
starting position \rightarrow ending position.
NB: $A[l, r]$ may be empty if no strings are prefixed by P .
- 2) ~~Pattern~~ String P is lexicographically located between $A[l-1]$ and $A[r]$.

EXAMPLE:

$$P = "te"$$



NB: There
5 STINGS
could be.
ANY WHERE
in memory.
POINTERS
(malloc)

$A[l, r]$

1-BASIC APPROACH's IDEA:

Use multiple use of P and $P\#$, where $\#$ is longer than any other character.

We search for:

- Startance Position I of PREFIX ARRAY $A[1, k]$ through BINARY SEARCH of P in D .
- Endance Position K of PREFIX ARRAY $A[1, k]$ through BINARY SEARCH of $P\#$ in D .

Two types of operations are supported.

① LISTANCE:

We need to perform:

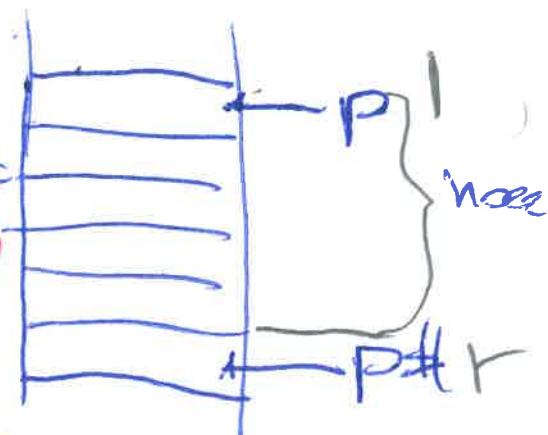
- ~~2 BINARY SEARCHES FOR TIME P and $P\# \Rightarrow O(P \cdot \log n)$~~

- LISTANCE of n_{oee} , to print the ~~all~~ strings between I and $K-1 \Rightarrow O(n_{oee})$ \Rightarrow We may have

$$\text{TIME} = O(P \cdot \log n + n) \quad \text{MISS} \text{ first thing omitted become } \text{ by PINTER INSTRUCTION.}$$

$$I/Os = O(P \cdot \log n + n) \quad \text{Non-CONTIGUOUS LISTING OF LINES}$$

Non-CONTIGUOUS LISTING OF LINES



② Containce:

We simply need to perform:

- 2 BINARY SEARCHES

- CONSTANT # OPERATIONS to find:

STRINGS needed by $P = t - 1$.

$$\text{TIME} = O(P \cdot \log n)$$

$$1/0s = O(P_B \cdot \log n)$$

SPACE by ① and ②: $O(M)$

No optimality!
↳ Poor scaling.

Worst # pointers.

2 - Containce Allocation of Strings.

NEW IDEAS: Store strings ~~contiguously~~
~~in lexicographical order~~ and more also
~~contiguously~~ on disk.

↳

Now: Pointers' contiguity reflect strings'
lexicographical contiguity.

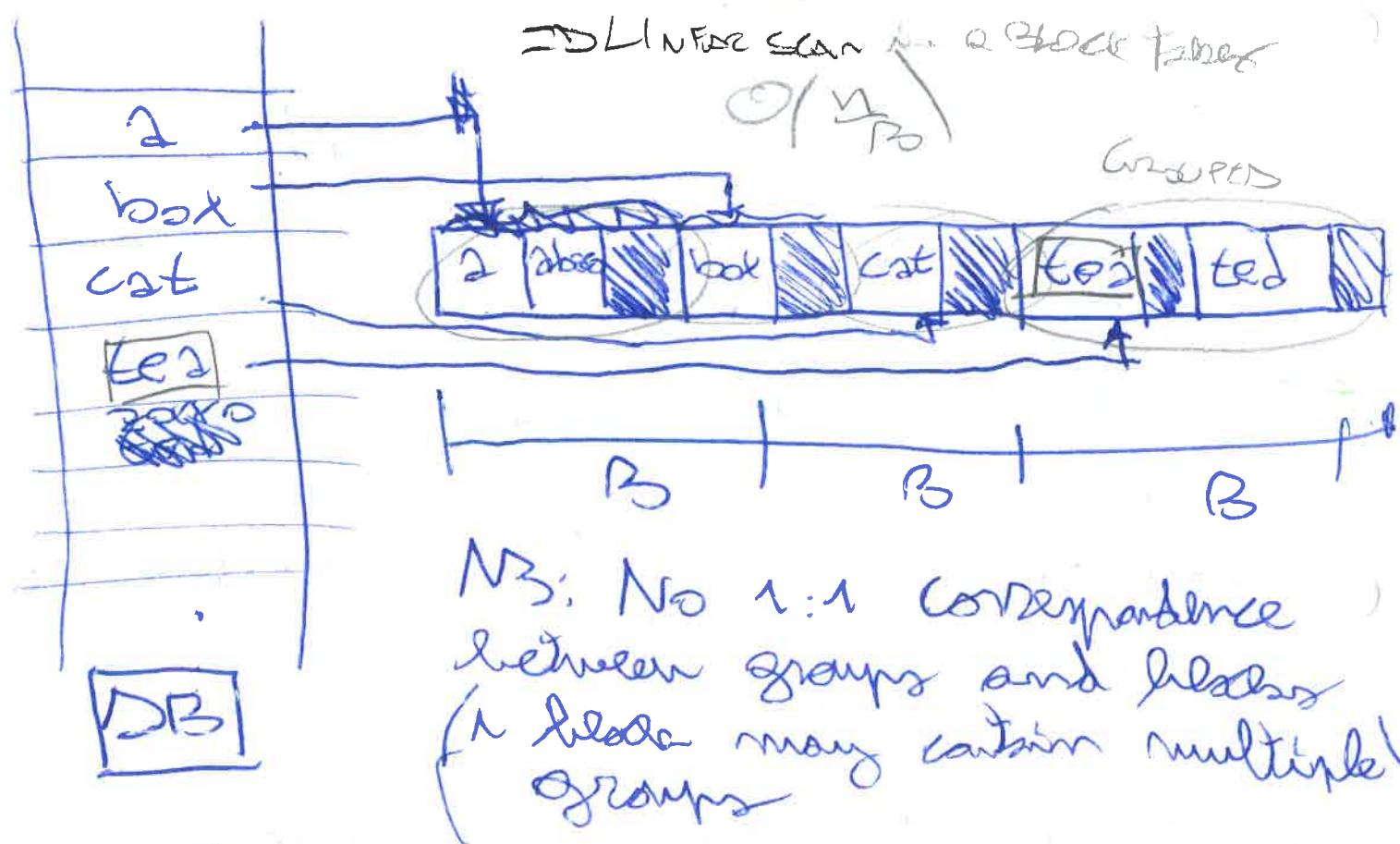
~~Advantages:~~

② When binary search is performed over
strings, it will benefit from a
SPEEDUP, as we will perform very
close groups in memory [very few 1/0s]

(D) We can apply some compression to contiguous strings, as they now share a PREFIX.

↳

Because of illustration ②, we notice that:
→ We are able to group contiguous strings into Blocks. In the DICTIONARY we just store a pointer to the first string in the block.



N.B.: No 1:1 correspondence between groups and blocks
{ 1 block may contain multiple groups}

→ We are hence able to deploy a 2-LEVEL SEARCH, which operates as follows:

Step 1: ~~in tree tables~~ Tries ΔP_B .

① Perform a DICTIONARY SEARCH of P over D_B to identify the lexicographic position of P (i.e. in which group it lies in).

STAGE

(3) Use the identified lexicographic position to find the blocks of strings where p lies in.

(3) Linear scan of the strings in the block to find either l or r .

NB: We repeat these (3) steps for both P and P^t .

~~Joining we can perform either LISANL or LISAND~~

HO cost for search phase of p 's last lexicographic position:

$$= \underbrace{\log_2 N}_{\text{PNNY}} - \underbrace{\log_2 B}_{\text{cost saving because over all pointers of blocks}} = \log_2 \left(\frac{N}{B} \right)$$

PNNY cost saving because over all pointers of blocks

We can perform either:

- LISANL:

$$\# \text{HOS} = O\left(\frac{P}{B} \cdot \log_2 \frac{N}{B} + \frac{\text{Mem}}{B}\right)$$

- LISAND

$$\# \text{HOS} = O\left(\frac{P}{B} \cdot \log_2 \frac{N}{B}\right)$$

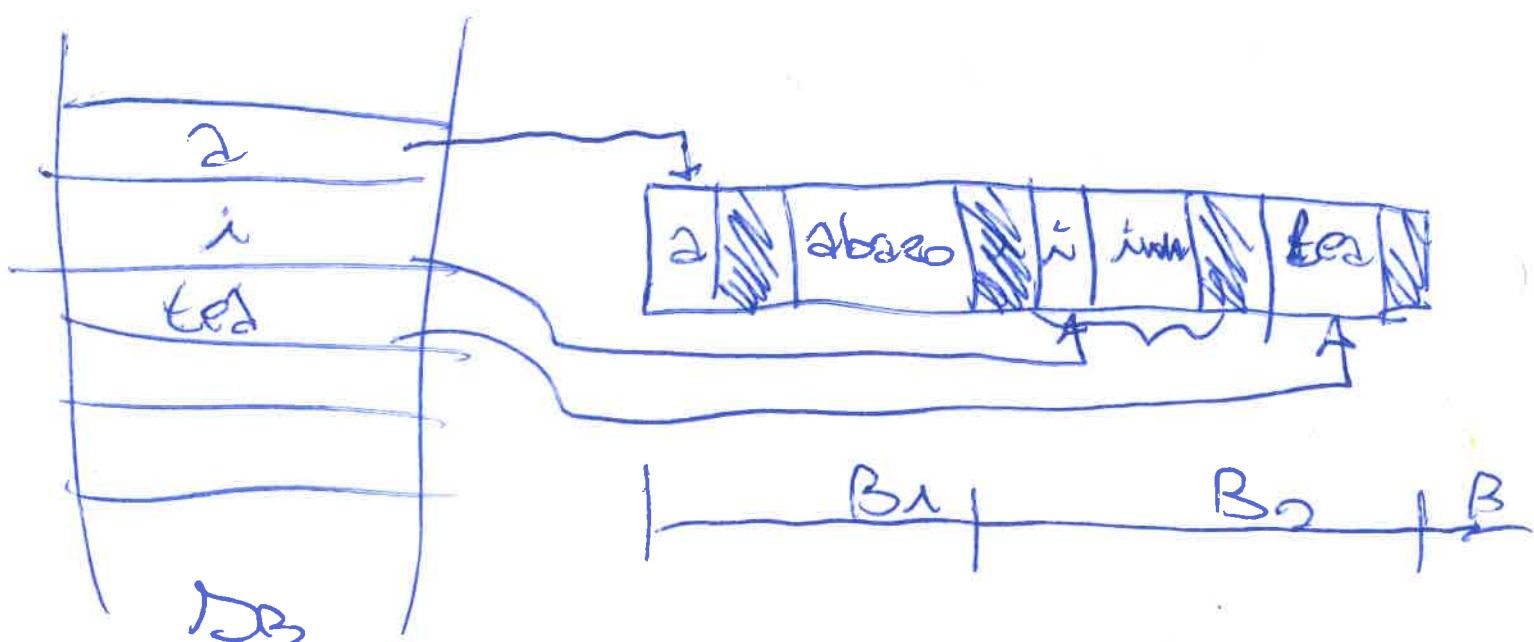
SPACE OCCUPANCY:

$$N + (1 + \text{rel. } n_B)$$

LISI all strings in lexicographical order of string

- $P \Rightarrow$ If things are very long (~~$\geq B$~~) ($\geq B$), Tiling technique has no advantages.
 \Rightarrow If things are short ($\leq B$), Tiling technique has advantages.

FUNCTIONING - VISUALIZAT.



~~(*)~~ $P = \text{"in"} \Rightarrow$ LISTING.

- ① Primary Search of P in DB
 \Rightarrow We get "i".
- ② We identify block of things, say (ii)
- ③ Linear scan of block to print "inh"

it takes O(1)

FRONT-CODING COMPRESSION (STANDARD)

Improves on the contiguous allocation system by ~~as well~~ lowering the space and time complexity further.

FDEXPLOR: Contiguous allocation of ~~strings~~ lexicographically sorted strings \Rightarrow These will likely share a prefix!

IDEC: We encode a string with as many possible characters from its preceding string (e.g.: LCP between strings S_{i-1} and S_i)
↳ Form of "D-compression algorithm".

$D = \langle \underline{\text{alcatraz}}, \langle \underline{\text{3,0d}}, \langle \underline{\text{alc0ne}}, \langle \underline{\text{anaclets}}, \langle \underline{\text{b}} \rangle \rangle \rangle$

$\left\langle \begin{array}{l} \langle \underline{\text{alcatraz}}, \langle \underline{\text{3,0d}}, \langle \underline{\text{alc0ne}}, \langle \underline{\text{anaclets}}, \langle \underline{\text{b}} \rangle \rangle \rangle \\ \downarrow I = \text{LCP}(S_{i-1}, S_i) \end{array} \right\rangle$
↳ Remaining characters.
(can be represented via $\text{LCP}(S_{i-1}, S_i)$)

ISSUE: There is a dependency of a string's encoding on the preceding strings.

Noise-case: To decode a string, we may need to go through all the preceding strings.

Ex D: $\langle 2 \rangle, \langle 22 \rangle, \langle 222 \rangle, \langle 2222 \rangle, \langle 22222 \rangle$
 $\langle 0, 2 \rangle, \langle 1, 2 \rangle, \langle 2, 2 \rangle, \langle 3, 2 \rangle, \langle 4, 2 \rangle$

→ WORST-CASE DECOMPOSITION:

↳ To decide $\langle \dots \rangle$, we need to go through all $|D|$ strings.

~~SOLUTION: FASTER~~

→ BETTER-CASE: String S_1 ~~obviously~~ has $LCP(S_{1:n}, S_1) \leq n \Rightarrow$ i.e. it is uncompressible \Rightarrow deciding it costs $O(n)$.

SOLUTION:

→ FRONT-CODING WITH BUCKETING.

NEW IDEA: We restart front-coding after the beginning of each block, where the first string of each block has left UNCOMPRESSED.

→ We compress other blocks individually and store the string length as well,
~~to repeat~~ to detect each string separately.

EXAMPLE:

$S = \langle \underline{\text{alex}}, \underline{\text{desk}}, \underline{\text{a}}, \underline{\text{z}} \rangle, \langle \underline{\text{alc}}, \underline{\text{o}}, \underline{\text{l}}, \underline{\text{co}}, \underline{\text{l}} \rangle, \langle \underline{\text{a}}, \underline{\text{kay}}, \underline{\text{a}}, \underline{\text{n}}, \underline{\text{o}} \rangle, \langle \underline{\text{a}}, \underline{\text{n}}, \underline{\text{a}}, \underline{\text{z}}, \underline{\text{e}}, \underline{\text{l}}, \underline{\text{o}}, \underline{\text{s}}, \underline{\text{z}} \rangle$

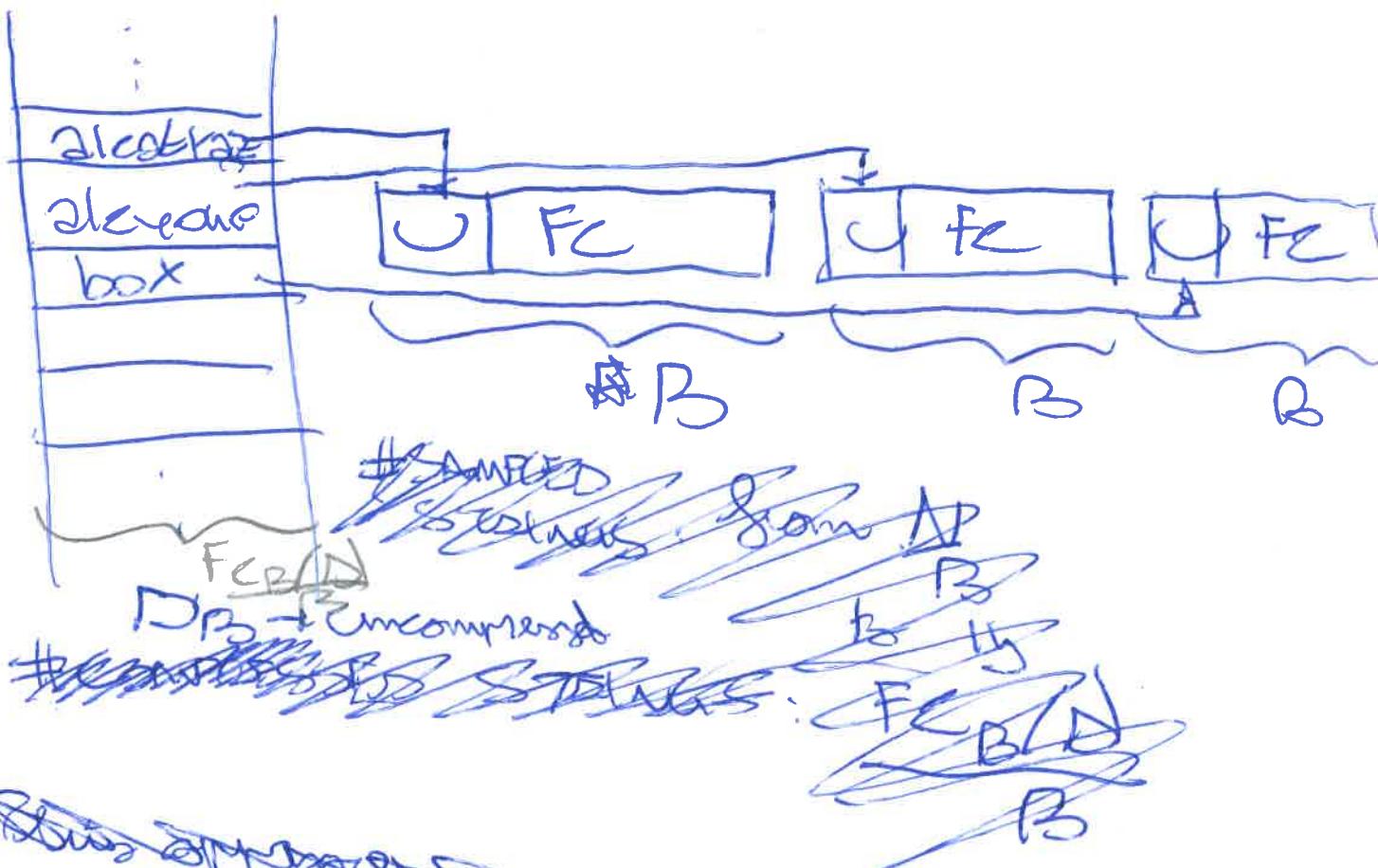


$\langle (B_1), \underline{\text{alc}}, \underline{\text{att}}, \underline{\text{a}}, \underline{\text{z}} \rangle \langle 6, 1, 0, 1 \rangle | \langle 7, 0, \underline{\text{alc}}, \underline{\text{ay}}, \underline{\text{a}}, \underline{\text{n}}, \underline{\text{o}} \rangle,$
↳ $\langle B_1, \underline{\text{alc}}, \underline{\text{ay}}, \underline{\text{a}}, \underline{\text{n}}, \underline{\text{o}} \rangle$
 $\langle B_1, \underline{\text{alc}}, \underline{\text{ay}}, \underline{\text{a}}, \underline{\text{n}}, \underline{\text{o}} \rangle$
 $\vdash LCP(S_{1:n}, S_1)$

→ IDEA: Use the compressed strings ~~and~~ ~~addresses~~ (first in the blocks) to allow retrievals.

The Block, Worst-case: Need to decompress all the preceding summary blocks & begining

NEW SITUATION VISUALIZED:



2 Phases are now:

Phase 1) We perform PARENARY SEARCH just over the UNCOMPRESSED STRINGS to identify

Phase 2) As larger block is compressed
intermittently, no shift is ever more than the block
itself to search.

(STRESS)

SAMPLED STRINGS after scan
NB: TOTAL SPACE till FCB(D) (NO FC with)
Blocks

$\text{FC}(B)$ = Space required by FCF to go
all the dictionary strings.

$\text{FCB}(D)$
B

For the
INDEXING
& DROPPING
UNCOMPR.

#I/O for Prefix Search = $O\left(\frac{P}{B} \cdot \log \frac{FCB(D)}{B}\right)$
 with FC - with Bucketing

STILL, we can still have some "worst case" ^{INDEXING} ~~Worst Case~~ worst case ~~making~~ $O(B^2)$ block access.

Ave. # Blocks needed	Worst-case # Blocks needed
"	"
$O(d)$	$O(B^2)$

$$DBS = \langle 2 \rangle, \langle 22 \rangle, \langle 222 \rangle, \langle 2222 \rangle, \dots$$

$$\#DBS = \langle 0, 2 \rangle, \langle 1, 2 \rangle, \langle 2, 2 \rangle, \langle 3, 2 \rangle, \dots$$

$$\#DBS = 0, 1, 2, \dots$$

, $K-1, K$

$$= \sum_{n=1}^K n = \Theta(K^2)$$

$$\text{Because: } \sum_{i=1}^n i = \frac{n(n+1)}{2}$$

GAUSS

$$\# \text{Worst-case shared characters} = K^2$$

- LISTING:

$$\# I/Os = O\left(\frac{P}{B} \cdot \log \frac{FCB(D)}{B} + \frac{FCB(D)}{B}\right)$$

- COUNTING:

LISTING from a Block!

$$\# I/Os = O\left(\frac{P}{B} \cdot \log \frac{FCB(D)}{B}\right)$$

LOCALITY-PRESERVING FRONT CODING

Elegant variation of front-coding, providing a trade-off between space occupancy and TIME to decode a string.

~~TO OPTIMIZE~~ wrt. #1/bs.

IDEA: A string S is front-coded only if its decoding time is proportional to its length $|S|$. Otherwise, it is written (uncompressed).

IDEAS: We compress a string S only if its decoding is **OPTIMAL**, otherwise we write it uncompressed (i.e. we copy it).

ALGORITHM: Δ STEPS handles block-level threading

if $d(S) \leq C \cdot |S|$ // SYSTEM PARAMETERS
copy S // WRITE UNCOMPRESSED
else:
 $FC(S)$

} MAX. as many steps as $C \cdot |S|$

OPTIMAL DECOMPRESSION: $\forall S \in D$:

#1/blocks in Decompression =

$$O\left(\frac{|S|}{B \cdot E}\right)$$

TOTAL SPACE = $O(1+E) \cdot FC(D)$

$$\epsilon = \frac{2}{L-2}$$

→ SPACE ~~depends~~
of compressed string
is also proportional to ϵ .

2) DECOMPRESSIVE hence costs:

- If S is UNCOMPRESSED

$$\# \text{1/OS} = O(\frac{|S|}{P_0})$$

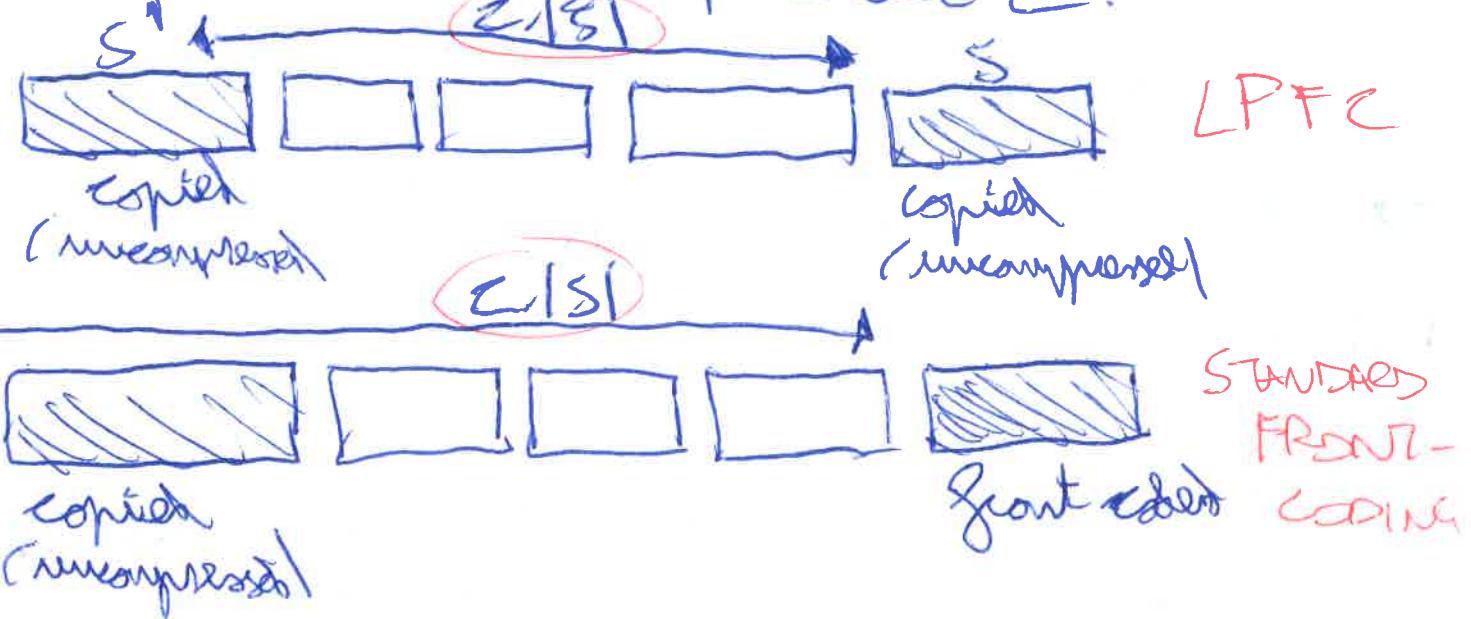
- If S is FRONT-CODED.

$$\# \text{1/OS} = O(\frac{C|S|}{P_0})$$

~~Proof: We want to consider two cases for copied strings.~~

~~- Case 1: #FC-extensions $\leq \frac{|S|}{2}$~~

We want to show that strings left uncompressed by LPFC (and were instead compressed by the classic front coding), have a length controllable through parameter C .



CLASS FOR RAGHDAF

CLASS EXERCISE - FRONT-CODING

$S = \{ \text{aabbb}, \text{aaabb}, \text{bbb}, \text{abbbbC}, \text{Abbbbd}, \text{abbbE}, \text{abc} \}$

FC-COMPRESSION (standard)

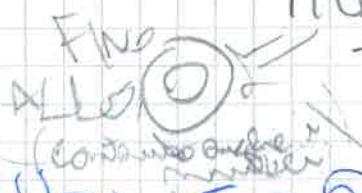
$\langle 0, \text{aabbb} \rangle; \langle 2, \text{bbb} \rangle; \langle 1, \text{bbbc} \rangle; \langle 4, \text{d} \rangle;$
 $\langle 4, \text{e} \rangle; \langle 2, \text{c} \rangle$

LPC with $C=1$

Recall the compression algorithm.

If $D(S) > C \cdot |S|$ // It stops you are allowed to go back by
 copying S
 else.

FC(S)



// CONSTRUCT ONE LEVEL

FOR i = 1 TO n

 CREATE

 PREFIX

 NODE

 # UNSET

 GIVEN

 TO

 WORK

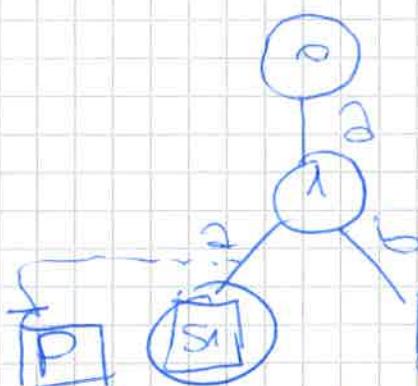
 OF

 PROGRAM

$\langle 0, \text{aabbb} \rangle; \langle 2, \text{bbb} \rangle; \langle 0, \text{abbbC} \rangle;$
 $\langle 2, \text{abbbd} \rangle; \langle 0, \text{abbbE} \rangle$

→ Build parallel tree on two strings

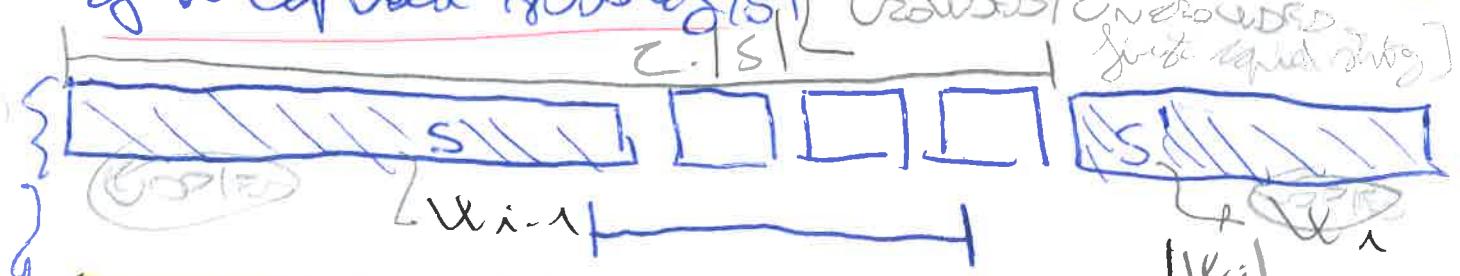
$S_1 = "aabbb"$ $S_2 = "abbbd"$
 Index $P = "aabbbd"$



(1) S_1 is reached. $\text{LCP}(\text{aabbb}, \text{abbbd}) = 22$
 $\Rightarrow l = 3$

Proof: That the total # characters in an uncrowded / crowded string is proportional to the length of its first copied string (uncrowded one)

\Rightarrow We consider two cases of copied strings, based on the ~~#FC-CHARACTERS~~, that lie between two consecutive occurrences of a copied string (S_i)



$$\text{CROWDED} \Rightarrow \# \text{FC-CHARACTERS} \leq \frac{C. | S_i |}{2}$$

(Following) before it



$$\text{UNCROWDED} \Rightarrow \# \text{FC-CHARACTERS} \geq \frac{C. | S_i |}{2}$$

\Rightarrow We would like to BEND the total length of copied strings into chains.

the note [uncrowded + crowded part] to be considered adding copied strings into chains.

$$\text{Chain} = [x_1, x_2, \dots, x_x]$$

(uncrowded copied string) (crowded copied string)

Take any copied string we have.
uncrowded x_i

$$x_i \geq \frac{x_{i-1} + x_i + x_i}{2} \Rightarrow \text{crowded}$$

$$|v_{i-1}| \geq \frac{z|v_i|}{2}$$

$$\Rightarrow |v_i| \leq \frac{2}{z} |v_{i-1}|$$

Let's consider the total length of a chain:

$$|v_1| + |v_2| + \dots + |v_x|$$

$$\sum_{i=1}^x |v_i| = |v_1| + \sum_{i=2}^x |v_i| \leq$$

$$\Rightarrow \sum_{i=1}^x |v_i| \leq |v_1| + \underbrace{\sum_{i=2}^x \left(\frac{2}{z}\right) |v_{i-1}|}$$

We can observe that:

$$|v_i| < \frac{2}{z} |v_{i-1}| < \left(\frac{2}{z}\right)^2 |v_{i-2}| < \left(\frac{2}{z}\right)^3 |v_{i-3}|$$

And therefore:

$$|v_i| < \left(\frac{2}{z}\right)^{i-1} |v_1|$$

And so, considering the total length:

$$\sum_{i=1}^x |v_i| \leq |v_1| + \underbrace{\sum_{i=2}^x \left(\frac{2}{z}\right)^{i-1} |v_1|}_{\left(\frac{2}{z}\right)^0 |v_1| + \sum_{i=2}^x \left(\frac{2}{z}\right)^{i-1} |v_1|}$$

$$\left(\frac{2}{z}\right)^0 |v_1| + \sum_{i=2}^x \left(\frac{2}{z}\right)^{i-1} |v_1|$$

$$\sum_{i=0}^x \left(\frac{2}{z}\right)^i |v_1|$$

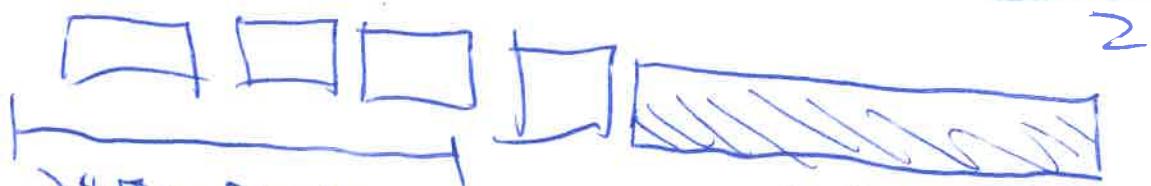
$$\Rightarrow \sum_{i=1}^x |v_i| \leq \sum_{i=0}^x \left(\frac{c}{2}\right)^i |v_i|$$

$$\sum_{i=1}^x |v_i| \leq |v_1| \cdot \frac{c}{c-2}$$

$$\frac{1}{1-\frac{c}{2}} = \frac{1}{\frac{c-2}{2}} = \frac{2}{c-2}$$

We have hence shown that every chain is bounded by $|v_1|$ and a constant!

\Rightarrow Since v_1 is UNBOUNDED, it will be preceded by at least $\frac{c|v_1|}{2}$ FC-characters



#FC-chars
i.e., we know that $v_1 \rightarrow$ UNBOUNDED v_1

$$\# \text{FC-chars} \geq \frac{c|v_1|}{2} \geq \frac{c}{2} \cdot \text{FC}(D)$$

And the # FC-characters before v_1 is bounded by $\text{FC}(D)$.

$$\text{FC}(D) \Rightarrow \text{FC}(v_1) \geq \frac{c|v_1|}{2}$$

\Rightarrow The length of the unbounded string v_1 is hence:

$$|v_1| < \frac{2}{c} \text{FC}(v_1)$$

~~#FC(v1) <= 2/2 * FC(v1)~~

* We have some test increased storage capacity.

$$|W_i| < \frac{2}{c} FC(D) < \frac{2}{c} FC(D)$$

And test:

$$\sum_{i=1}^x |W_i| < \text{test. } \frac{2c}{c-2}$$

PLC ~~values~~ $\sum_{i=1}^x |W_i| < \cancel{\frac{2}{c} FC(D)} \cdot \cancel{\frac{2c}{c-2}}$

$$\sum_{i=1}^x |W_i| < \frac{2}{c-2} FC(D)$$

$$\Rightarrow \Sigma = \frac{2}{c-2} O.E.D.$$

ADVANTAGES ~~of~~ PFC.

Compressed storage, acting as space booster.

We are able to uncompress anything at any point in time, and we do ~~not~~ ^{"parallel"} cases.

→ No slowdown on complexity!

Optimal I/Os & time!

INTERPOLATION SEARCH

IDEA: Combine **BINARY SEARCH** with **INTERPOLATION** to achieve a very fast search strategy.

INPUT

Sorted array $X[1 \dots m]$, array size

B_1

B_3

B_5

B_7

B_9

B_{11}

B_{13}

B_{15}

B_{17}

B_{19}

B_{21}

$X =$	1	2	3	8	9	17	19	20	28	30	32	36
	1	2	3	4	5	6	7	8	9	10	11	12

Buckets' creation

X is sub-divided into ~~BUCKETS~~ $\frac{X_m}{b}$ of size b .

$$b = \frac{X_m - X_1 + 1}{m} = \frac{X_{\max} - X_{\min} + 1}{m}$$

because sorted array

1 block (and length element)

EXAMPLE:

$$b = \frac{36 - 1 + 1}{12} = 3 \Rightarrow \text{BUCKET }$$

$$\# \text{BUCKETS} = \frac{X_m}{b} = \frac{36}{3} = 12$$

BUCKETS: $[19-21] [22-24] [25-27] [28-30]$
 $B_7 \quad B_8 \quad B_9 \quad B_{10}$

$[1-3] [4-6] [7-9] [10-12] [13-15] [16-18]$
 $B_1 \quad B_2 \quad B_3 \quad B_4 \quad B_5 \quad B_6$

Now also consider an array I , containing the starting and ending position of every bucket.



Ques: How to search how to jump to a certain bucket in $O(1)$ time.

$I \rightarrow$ starting & ending pos. of every bucket

(Start pos, End pos, Length)

PSEUDOCODE:

ML Search(y):

① $j = \lfloor \frac{y - x_1}{b} \rfloor + 1$ // j-th BUCKET
for

② Binary Search ~~for~~ y in the portion of X delimited by $I[j]$

① Find Bucket By

② Binary search for y over B_j .

TIME: $O(\text{FIND BUCKET } B_j) + \text{BINARY SEARCH over } B_j$
(PRELIMINARY)
2. Really ESTIMATE

$$O(1 + \log_2 b) = O(\log_2 b)$$

SPACE: $O(m)$

My pretty inaccurate estimate of the time complexity.

We need to take the time complexity of the buckets & elements therein.

EXERCISE (IN CLASS)

INTERPOLATION SEARCH EXERCISE.

1, 3, 4, 9, 16, 18, 20, 21, 28, 30, 32, 36

- (1) Partition these into BUCKETS, where the BUCKET SIZE is:

$$b = \frac{X_m - X_1 + 1}{m} = \frac{36 - 1 + 1}{12} = 3$$

$$\# \text{BUCKETS} = \frac{X_m}{b} = \frac{36}{3} = 12, \text{ with } 3 \text{ elements each.}$$

BUCKETS' REPRESENTATION → 2 WAYS:

- (1) <Starting position, bucket length>
<element>

- (2) <Starting position, ending position>
<element>

→ Use (2) for REPRESENTATION (1).

- (2) Fill-in values into the corresponding bucket.

1 1, 3	4 4	7 9	10 10	13 9
16 18	19 20, 21	22 23, 24	25 26	28 29, 30
31 32	34 35	37 38	40 41	43 44

③ Look for ITEM : (2) (for instance).

→ We retrieve bucket (1), then perform a BINARY SEARCH into this bucket over [20 | 21] → (1 time!)
no item retrieved!

GOAL: We want to better estimate the runtime of INTERPOLATION SEARCH.

TIME for INT. Search = $O(m + \log_2 b) = O(\log_2 b)$

Find
Bucket by
Binary Search
over bucket by.
↳ DOMINATES the
cost.

THEOREM:

(INTERPOLATION)

Binary search cost

in dictionary of size m .

$O(\log \Delta)$ in the worst-case, where

$$\Delta = \frac{\text{MAX. GAP between } x_i \text{ and } x_{i-1}}{\text{MIN. GAP between } x_i \text{ and } x_{i-1}}$$

EXTRA SPACE = $O(m)$, (OPTIMAL ASYMPTOTIC)

PROOF: CORRECTNESS is immediate.

TIME COMPLEXITY: The maximum of a series of integers is at least as large as their average (GAP).

- MAX. GAP \geq AVG. GAP

$$\frac{\max(x_i - x_{i-1})}{m} \geq \frac{1}{m} \sum_{i=2}^{m-1} x_i - x_{i-1} = \frac{1}{m-1}$$

MAX. GAP \geq AVG. GAP

$$\max_{i=2,..,m} X_i - X_{i-1} \geq \sum_{i=2}^m \frac{X_i - X_{i-1}}{m-1} = \frac{X_m - X_1}{m-1} + 1$$

TELESCOPIC SERIES

MAX.
BIN
SIZE

$$\Rightarrow \max_{i=2,..,m} X_i - X_{i-1} \geq \frac{X_m - X_1 + 1}{m} = b$$

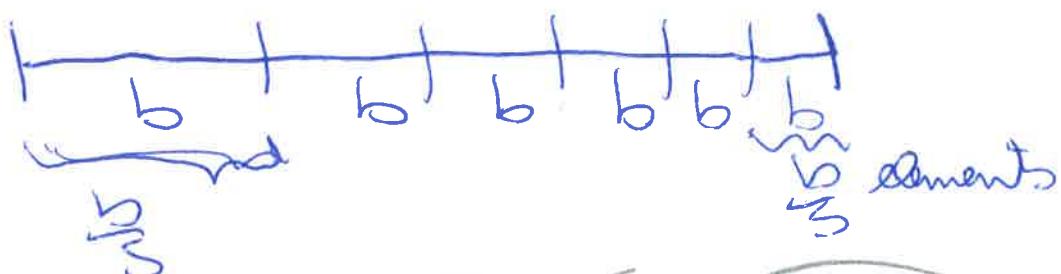
UPPER BOUND

- Also, the ~~Max. # INVERSIONS~~ test can belong to any bin b can be upper bound.

\Rightarrow Integers of X are spaced apart by:

$$s = \min_{i=2,..,m} (X_i - X_{i-1}) \text{ CNTS per BIN.}$$

\Rightarrow Every bin of size b has at most $\frac{b}{s}$ elements.



$$\Rightarrow |B| \leq \sum \frac{b}{s}$$

MAX. BUCKET SIZE

\Rightarrow Let's put together the lower & upper bounds in the definition of Δ .

$$D = \frac{\text{Max. GAP between } x_i \text{ and } x_{i-1}}{\text{Min. GAP between } x_i \text{ and } x_{i-1}}$$

~~is it true? And we have found that:~~

$$\text{as } |B_i| \leq \frac{b}{s} \leq \frac{\max_{i=2,..,m} (x_i - x_{i-1})}{\min_{i=2,..,m} (x_i - x_{i-1})} = D$$

$$|B_i| \leq \frac{b}{s} \leq D$$

\Rightarrow The Algorithm's complexity hence depends on the SPARSENESS of values, not on their MAX/MIN value.

• The more sparse items are, the faster the Algo. Runs [few items in ~~one~~ bin] \Rightarrow DBLMAX to D.

WORST-CASE SEARCH TIME: If whole ~~area~~ ends up in a single BIN, the worst-case search time is hence:

$$\mathcal{O}(\log \min(D, m))$$

[No worse than BINARY SEARCH] Same as performing BINARY SEARCH

\Rightarrow We are more interested in the AVERAGE SEARCH TIME.

AVERAGE SEARCH TIME.

The time complexity is $O(\log \log m)$,
if ~~M~~ the integers are drawn uniformly at random (i.e. uniform distribution).

Proof: We want to estimate D with high probability.

(NB: $\log \log m \ll \log m$)
one order of magnitude

\Rightarrow Assume integers ~~integers~~ are drawn uniformly from $[0, U-1]$.



Divide the partition the range into ~~ranges~~ ^{to the} ~~where the sought values lie~~

k ranges \Rightarrow Each range has $\frac{U}{k}$ elements,
where $k = \frac{m}{2 \cdot \log m}$
 $= \# \text{ranges.}$

We have k ranges.

$\Pr[\text{one integer belongs to a range}] = \frac{1}{k}$
out of k existing ranges

$P\{$ one specific range does not contain one integer } = $1 - \frac{1}{m}$

$P\{$ one specific range does not contain any integer } = $(1 - \frac{1}{m})^t$

We know that

$$t = \frac{m}{2 \cdot \log m}$$

$$\Rightarrow P\{ \text{one specific range does not contain any integer} \} = (1 - \frac{2 \log m}{m})^{2 \log_2 m}$$

$$= \left[(1 - \frac{2 \log m}{m})^{\frac{m}{2 \log_2 m}} \right]^{2 \log_2 m}$$

$$= e^{-2 \log m}$$

$$\Rightarrow O(e^{-2 \log m}) = O(\cancel{e^{-2 \log m^2}})$$

$$= O\left(\frac{1}{m^2}\right)$$

$\Rightarrow P\{ \text{at least one range remains} \} \leq O\left(\frac{1}{m}\right)$

} (or, ~~empty~~ range contains at least one element)

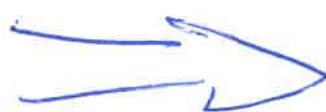
If this occurs with high probability, then the **max. distance** between adjacent integers must be smaller than twice the range's length.

$$\max_i x_i - x_{i-1} \leq 2 \cdot \frac{U}{T} = \frac{2 \cdot U \cdot 2 \cdot \log m}{m}$$

$$t = \frac{m}{2 \cdot \log m}$$

$$\Rightarrow \max_i x_i - x_{i-1} = O\left(\frac{U \cdot \log m}{m}\right)$$

Let's now take $t' = O(m \cdot \log m)$ ranges.



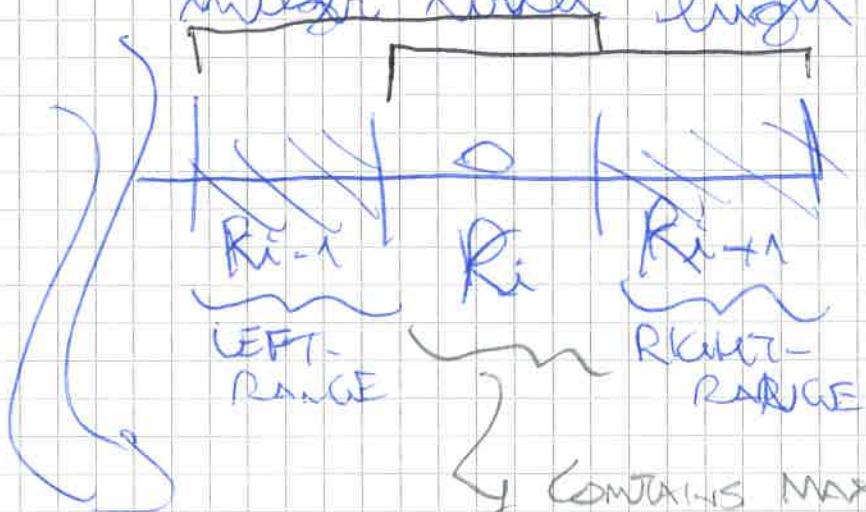
We now want to find

$$\min_i x_i - x_{i-1}$$

\Rightarrow Now take $f' = \Theta(m \cdot \log m)$ ranges

\Rightarrow (less than m in total)

We can prove that every adjacent pair of ranges contains at most one integer with high probability



R_i = i th range

containing ~~max~~ 1 integer with high probability!

IF a range R_i contains integer, its two adjacent ranges (LEFT & RIGHT) are empty with high probability.

$$\Rightarrow \min_{i=2}^m x_i - x_{i-1} \geq \frac{f'}{f'} = \Theta(\frac{1}{m \cdot \log m})$$

~~minimum distance between an adjacent pair of integers~~

Let's now take D :

$$D = \frac{\max_i x_i - x_{i-1}}{\min_i x_i - x_{i-1}}$$

$$\leq \frac{\Theta(\log m)}{\Theta(\frac{1}{\log m})} = \Theta(\log^2 m)$$

INTERPOLATION SEARCH'S

TIME COMPLEXITY $D = O(\log^2 m)$ |
 even much prob |

#110s: $\mathcal{O}(\frac{P}{B}, \log \log \frac{N}{P})$

→ EXP. reduction in search time performance! (more binary search)

N.B.: If strings are NOT uniformly distributed

Randomly

→ We permute them to get a UNIFORM DISTRIBUTION

COMPACTED TREE:

Compressed trees consist of EFFICIENT data structures for SEARCHING (strings).

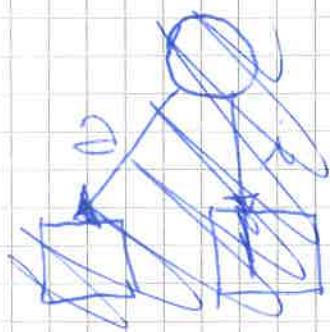
TREE → Index compressed strings in internal memory

Split up in first place & modify search from $\mathcal{O}(\log(\frac{N}{P}))$ to
~~(Binary Search)~~ $\mathcal{O}(P)$ TIME.

Independent from dictionary size, as in RAM we can manage and address memory cells of $\mathcal{O}(\log N)$ bits

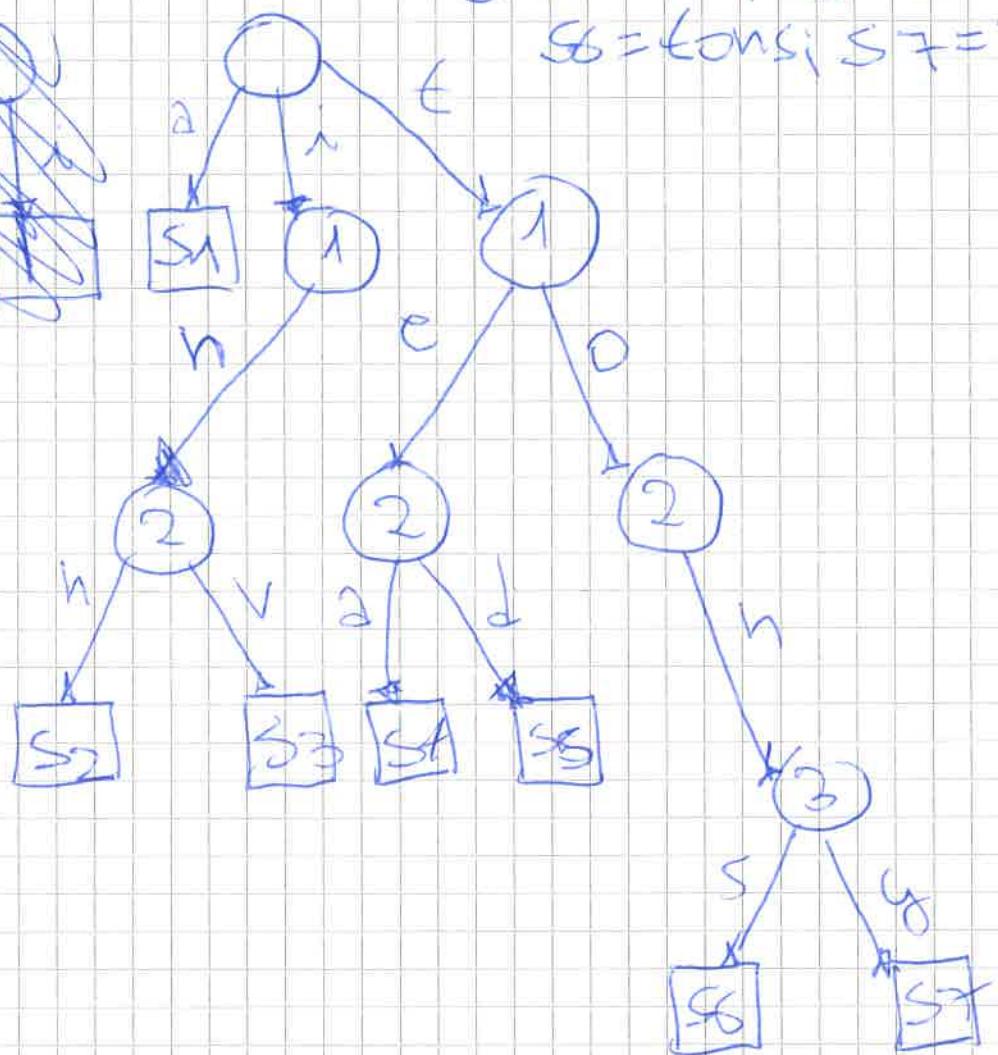
— Also save space! No unnecessary paths, as there are compacted!

UNCOMPACTED
TRIE
(not):

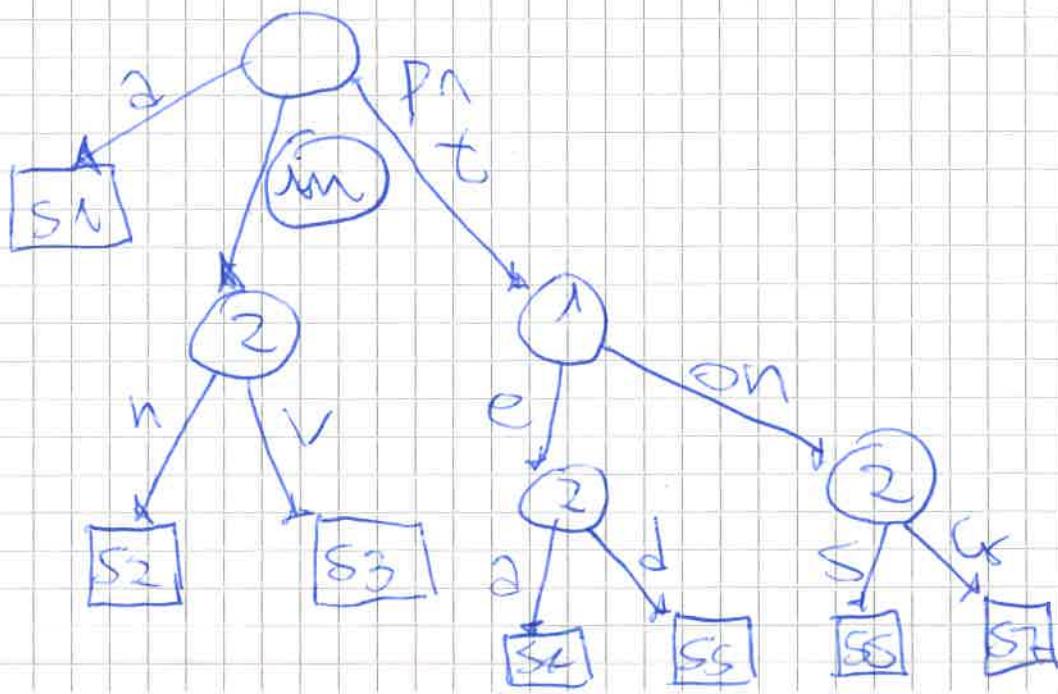


Stringo:

$S_1 = \text{a}$; $S_2 = \text{inn}$; $S_3 = \text{inv}$;
 $S_4 = \text{ea}$; $S_5 = \text{ed}$;
 $S_6 = \text{kons}$; $S_7 = \text{tony}$



Corresponding COMPACTED TRIE.



COMPACTED TREES FEATURES.

LEAVES = n

INTERNAL NODES < n

NODES + LEAVES < 2n

SPACE = $O(n) + O(N)$ ($\# \text{NODES} = O(n)$)
 ~~$O(n) + O(N) \rightarrow \text{Space is better}$~~

PREF(k)-SEARCH TIME = $O(p + ncc)$

IPOS: $O(p + ncc)$

↳ Leaves & Internal Nodes
 Contiguously & sorted alphabetically
 only on links }

SEARCHES
 required by
 pattern P.

EDGES' REPRESENTATION.

Edges are simply represented as pointers
 to (a part of) characters in leaves' strings.

Ex. P1 = "in" = $\langle S_2, 1, 2 \rangle$ → starting index of the substring

Ex2: P2 = "on"
 ↳ Starting index of the substring

$\langle S_2, 1, 2 \rangle$ $\times 12$ bytes
 + 4 bytes for string bytes

Hence, the space required by edges is:

$n_2 \cdot \# \text{EDGES} = O(n)$

(Proportional to the ~~# EDGES~~
 # STINGS)

\rightarrow SPACE is hence proportional to the #STRINGS that can be stored in internal memory (otherwise, linear dependence of #1/0s on pattern length p)



#Strings that can fit into the INDEX.

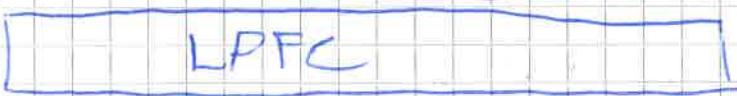
Ex. 2 GB Ram \Rightarrow 2,000,000 strings can fit in internal memory & be properly indexed.

M



\Rightarrow ENCODED TRIE

→ Strings stored in M.

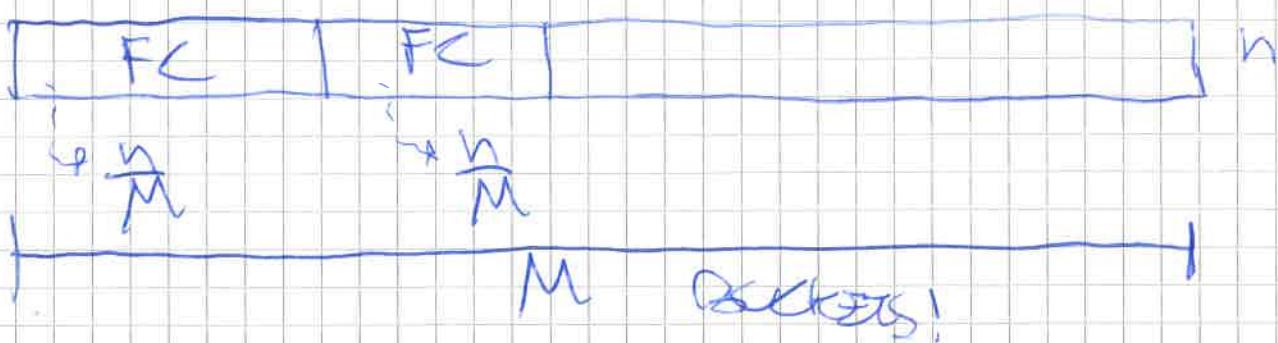


ISSUE: What if some strings do not fit in the internal memory (M)? (SAME APPROACH AS FOR COMPACTED TRIES)

(i.e.: If $M \rightarrow M'$)

→ STRINGS also need to be stored on disk, "SAMPLE".

See many make use of BUCKETS.



PATRICKA TRIES

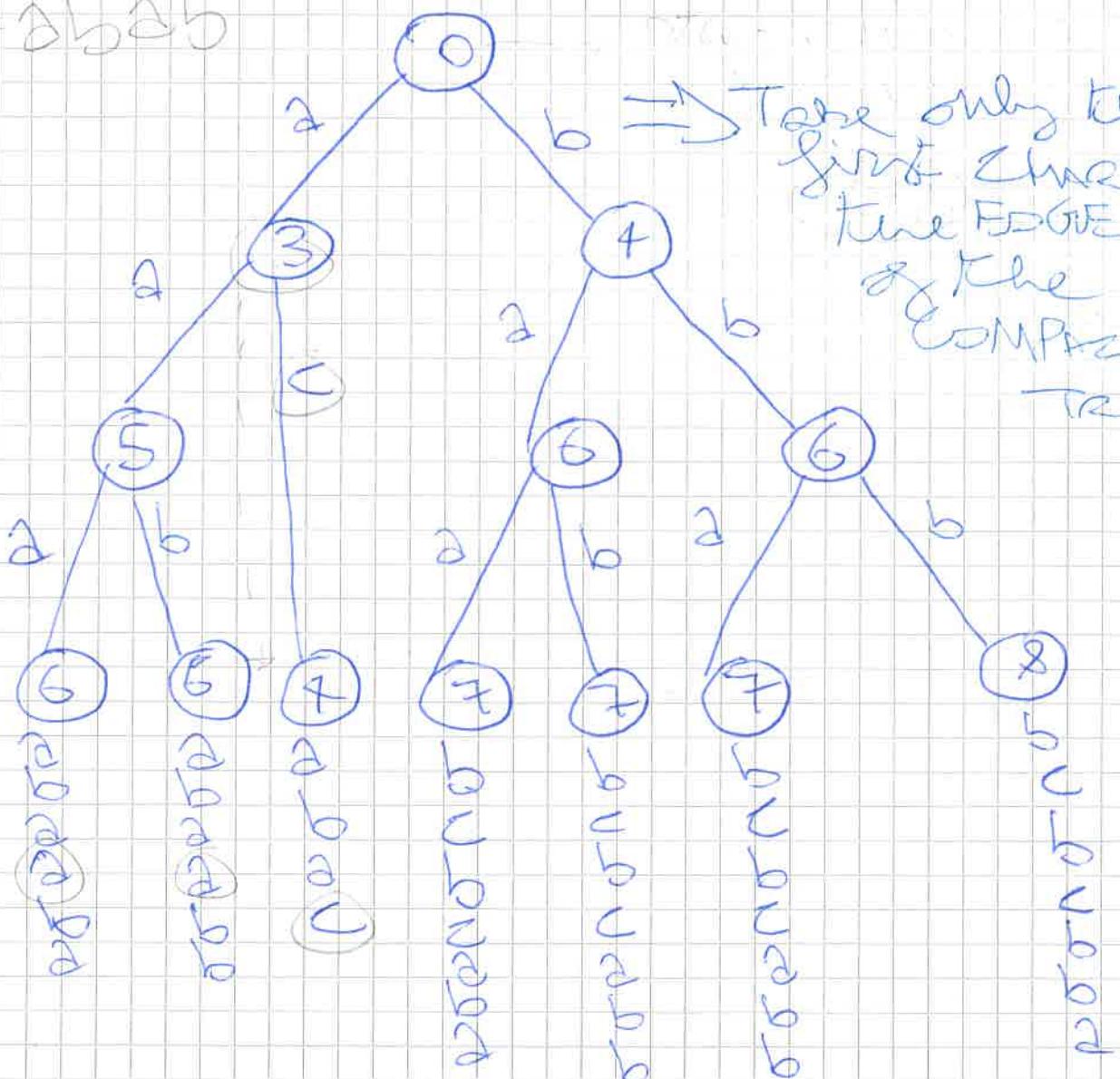
Data structure that, after being built, allows for search in $O(1)$ TIME & $1/258$ of a pattern among a sorted sequence of strings.

ADVANTAGE: Access one single string in $O(1)$ instead of $O(|P|)$.

PATRICKA TRIE

$P = abab$

PATRICKA TRIE:



COMPACTED TRIE

SEARCH PROCESS IN A PATRICK TRIE.

INPUT: pattern P , Patricia TRIE just 2
SINK

- 1) In $O(N)$ time, trace $S \rightarrow$ DOWNWARD path to a leaf, say l , until no branching is possible any longer.
 [Keep going down as long as characters match along the ~~other~~ edges].
 If ~~tie~~ \Rightarrow tie among leaves, pick one Randomly.

- 2) In $O(|P|)$ time, compute:

$$L = \text{LCP}(P, S)$$

\Rightarrow We will hence know that:
 $S[L+1]$ is the ~~mismatch~~ mismatching character.

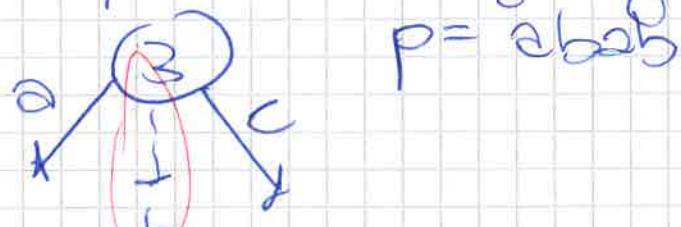
$$P[L+1] \neq S[L+1]$$

N.B.: Leaf l stores the string having its LCP of P and S .

- 3) Perform an UPWARD TRAVERSAL from leaf l up to an EDGE where the mismatching character $S[L+1]$ lies.

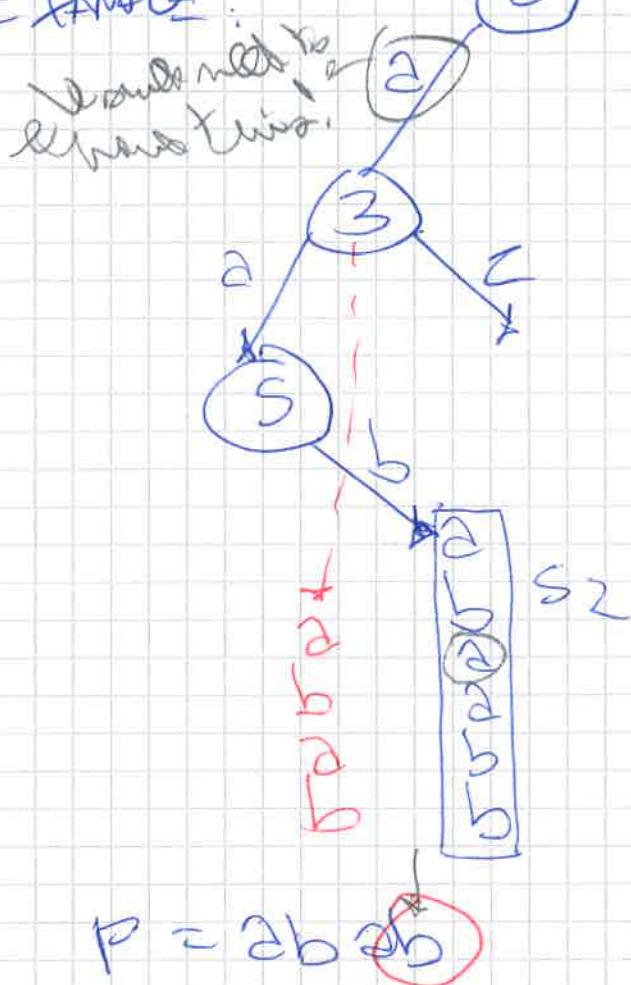
(i.e.: go up and find $S[L+1]$)

(4) Identify the correct branching position for pat the edge found.



$$p = abab$$

EXAMPLE:



$$P = ab\textcolor{red}{ab}$$

(1) Downward pass: Reach $S[3]$

$$(2) L = \text{LCP}(P/S_2)$$

$$= \text{LCP}(\underline{ababbb}, \underline{abab}) = 3$$

(3) Go up and find mismatching character at $S[L+1]$ ~~two~~

(4) Plug in the ~~underwater~~ node (2).
Put the correct position.

COSTS in PATRICIA TREE:

The blind search (downward pass) in a PATRICIA TREE takes $O(p)$ time and $O(p_B)$ #I/Os to compare things identified by blind search.

↳ Searching for p & $p\#$ insertion in the compacted TRE, we can find the differences between p and $p\#$.

$$\text{SPACE} = \underbrace{O(n)}_{\text{NODES}} + \underbrace{O(N)}_{\text{STRINGS}} \\ \text{& EDGES (no leaf)}$$

$$\text{SEARCH TIME} = \underbrace{O(p)}_{(\text{IN MEMORY})}$$

$$\# \text{I/Os for Search} = \underbrace{O(p_B)}_{\substack{\Rightarrow \text{Because of} \\ \text{LCP} \\ \text{computation} \\ \text{when reaching} \\ \text{the leaf}}}$$

USAGE of LPFC WITH PATRICIA TREE:

⇒ If we want to compare strings in English by the PATRICIA TREE, we proceed as follows:

- We put them ~~internal memory~~ in $\frac{\# \text{strings}}{1/8 n} > \text{Mts INTERNAL}$ memory the PATRICIA TREE of dictionary D.
- We store on disk the LPFC & the dictionary D.

\Rightarrow In this manner, traversing the PATRICKA TRIE takes $O(p)$ and No I/Os, as we fully operate in RAM.

\Rightarrow LCP computation takes:

$$O\left(\frac{L}{B} + \frac{P}{P_0}\right) \quad \# I/Os.$$

\Rightarrow Needs to decode strings from their LPFL-Representation.

~~WORD SEARCH~~

SPACE: $O(n)$

Patricia is embeddable in internal memory.

$$O((1+\epsilon)FC(D)) \text{ on } \underline{\text{disk}}$$

$$\Rightarrow \# I/Os = O\left(\frac{(1+\epsilon) \cdot FC(D)}{B}\right)$$

$\forall n = \sum M_i$ (i.e. strings of the PATRICKA TRIE do not fit in internal memory)

\Rightarrow We need to perform some "BUCKETING" and embed a sample of strings in the PATRICKA TRIE.



(Ex. First big & heavy Bucket embeddable in PATRICKA TRIE)

EXERCISE (IN CLASS)

a) Build a PATRICKA TRIE for:

$$S = \{ \underline{\text{0001100}}, \underline{\text{0001110}}, \underline{\text{0010}}, \underline{\text{101}}, \underline{\text{111}} \}$$

It is useful to build first the corresponding COMPACTED TRIE, and then the PATRICKA TRIE.

$$S = \{ \underline{\text{0001100}}, \underline{\text{0001110}}, \underline{\text{0010}}, \underline{\text{101}}, \underline{\text{111}} \}$$

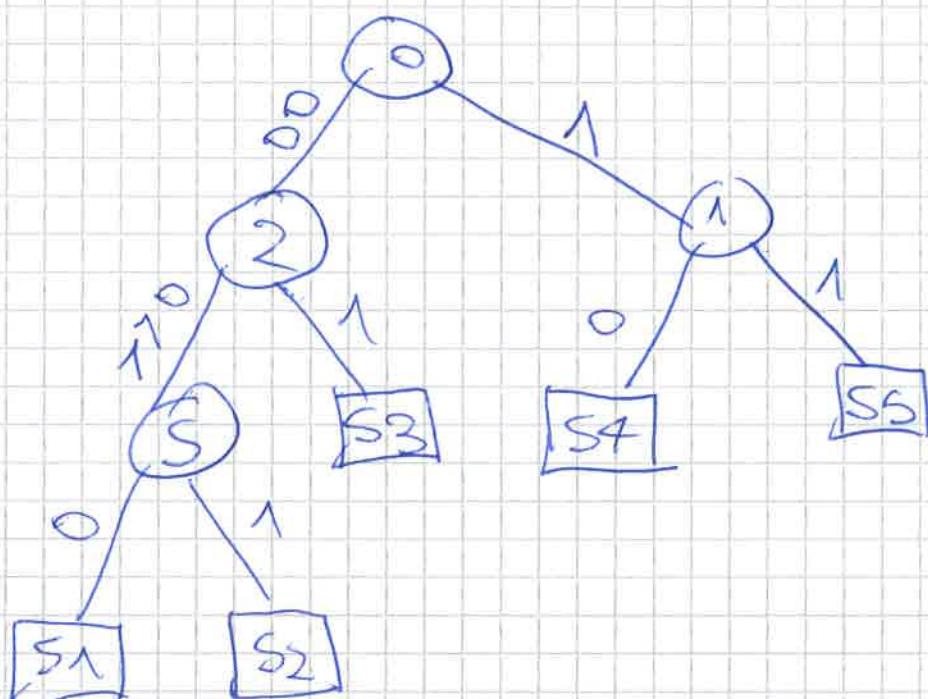
$\boxed{S1}$ $\boxed{S2}$ $\boxed{S3}$ ~~$\boxed{S4}$~~ $\boxed{S5}$

\Rightarrow Nodes have the # common characters

Edges have the common characters themselves

Leaves have the strings we want to store in the TRIE.

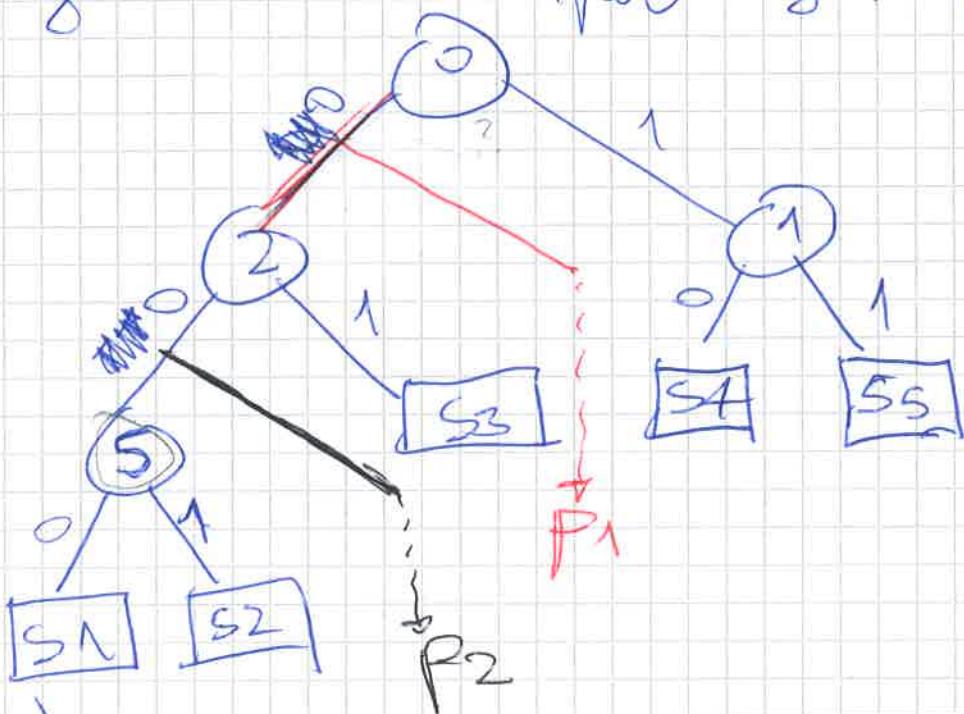
COMPACTED TRIE:



\Rightarrow The corresponding PATRICKA TRIE can be built simply by taking the first character of each edge from the COMPACTED TRIE.

PATRICIA TRIE:

It is the COMPACTED TRIE, but with the first character per slope only.



b) Find the LEXICOGRAPHIC Position for:

$$\textcircled{1} \quad P_1 = 0100\ 00$$

$$\textcircled{2} \quad P_2 = 0101$$

$$\textcircled{1} \quad P_n = 0100\ 00$$

\) STEPS: a) Downward Path

(Navigate PATRICIA TRIE ~~downward~~ as long as you find matching characters)

$\Rightarrow [S1]$ returned

$$\text{b) } LCP(P_n, S_n) = 1 = L$$

$$L = LCP(0001100, 010000)$$

c) Upward traversal up to the mismatching character ~~sign~~. $S[L+1]$

~~S1 = 0001100
P1 = 0101010~~

$$P1[2] < 0?$$

$L \leftarrow 1$

$1 < 0?$ No!

\Rightarrow Right by the slope!

2) $P2 = 010101$

- Downward path:

We move down to (S) \Rightarrow We can pick either (S_1) or (S_2) \Rightarrow We pick S_1

$$- L = \text{LCP}(P2, S_1)$$

$$= \text{LCP}(\underline{010101}, \underline{0001100}) = 1$$

~~Right~~

\Rightarrow Where do we put P_2 now?

\Rightarrow Go up to node 2

LCP in Linear Time - KASAI's ALGORITHM.

~~Suffix Array~~ ~~Linear~~ ~~Comparisons~~ $S = aabbzbaak$
 $\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \end{matrix}$

SUFFIX	SA	LCP	#LINEAR COMPARISONS	$O(n)$	$O(n^2)$
SUF(π)				ROUTE-FOLGE #COMPARISONS	
SA^{-1}					
$aabbzbaak$	1	3	4		
$aabbzbaak$	2	2	3		
aak	6	2	1		
$aabbzbaak$	3	1	1		
aak	7	1	1		
$aabbzbaak$	4	0	1		
K	8	0	1		
$zbaak$	5	-	-		

$SA^{-1} \Rightarrow$ Given an SA,
it returns its Rank.

$n = \# \text{SCRELY SUFFIXES}$
CLASSES.

SA^{-1} \downarrow LINEAR COMPARISONS

1	1	$SA=2 \{ aabbzbaak \} \Rightarrow \# \text{comps} = 4$
2	2	$\# LCP = 3$
3	4	$h=3$
4	5	$h=3$
5	8	$\Rightarrow \# \text{comps} = 1$
6	3	
7	5	
8	7	

$SA=4 \{ h=1 \}$
 aak
 $aabbzbaak \Rightarrow \# \text{comps} = 1$

$SA=5 \{ h=0 \}$
 K
 $zbaak \Rightarrow \# \text{comps} = 1$

$S_A = 6$

$\begin{cases} \text{a a a a z k} & \Rightarrow \# \text{cons.} = 3 \\ \text{a a k} & h = 2 \end{cases}$

$\begin{cases} S_A = 7 \\ h = 2 \uparrow \\ \text{a a a a a k} & \# \text{cons.} = 1 \\ \text{a k} \end{cases}$

$\begin{cases} S_A = 8 \\ h = 0 \\ \text{a z a a k} & \# \text{cons.} = 1 \\ k \end{cases}$

\Rightarrow IDE: Process by groups of 2, & check
LCP between Suff_i and Suff_{i+1}

⑩ SUBSTRING SEARCH.

We are presently interested in solving the FULL-TEXT SEARCH or SUBSTRING SEARCH problem.

INPUT: Text string $T[1, n]$

We want to:

- RETRIEVE (or count) all text positions where a query pattern $P[1, p]$ occurs as a substring of T .

Ex: $T = \text{trends}$
123456

$P = \text{re}$

OUTPUT: $T[s, e]$ or \textcircled{A}

"# occurrences of P in T .

SUFFIX-DEFINITION.

Consider a text T . Note that if suffix is a PREFIX of another suffix (abinga) then it ends at position n .

$\text{SUFF}_n = T[i, n]$ (i.e. Suffix i occurs at positions i, n)
The opposite of a PREFIX.
→ Occurs from i to n .

Ex: $T = \text{mississippi}$
123456789012

Suffix-6 = sippi\$

PREFIX of a SUFFIX:

$$\text{Is } P = T[i, i+p-1]$$

↳ Pattern $P[i, p]$ is a PREFIX of ~~the string~~

$$\text{SUFFIX} = T[i, n]$$

$$T = \text{"mississippi"} \\ \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \end{matrix}$$

$$P = "iss"$$

⇒ P occurs at pos. 4 and it prefixes
the string $T[4, 12]$
SUFFIX 4.

BRAVE-FORCE APPROACH FOR
SUBSTRINGS SEARCH:

⇒ Compare P against all possible
substrings of $T[i, n]$

⇒ Worst-case: $\mathcal{O}(n \cdot p)$

$$T = \text{mississippi\$} \\ \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \end{matrix}$$

$$P = "iss"$$

⇒ n comparisons, taking P each one.

- (1) mississippi; "iss"
- (2) ississippi; "iss"
- (3) ssissippi; "iss"
- (4) mississippi; "iss"
- (5) issippi; "iss"
- (6) ; "ississippi" and "mississippi" are returned.

SMARTER APPROACH - DATA STRUCTURES:

We consider some data structures to quickly retrieve & efficiently store SUFFIXES (i.e. SUFFIXES)

4. SUFFIX ARRAY (SA[])

It is an array of pointers to all text suffixes, ordered Lexicographically.

$SA[i]$ = ith-smallest text in lexicographical order.

SUFF
 $SA[1] \leftarrow \text{SUFF}$ $SA[2] \leftarrow \dots \leftarrow \text{SUFF}$ $SA[n]$
& Sorting by Lexicographical order.

SPACE = $O(n \cdot \log n) + O(n \cdot \log n)$

SPACE in terms of pointers
to elements of SA

Pointers to Suffixes
of SUFFIXES

Ex: T = mississippi\$
1 2 3 4 5 6 7 8 9 10 11 12

SORTED SUFFIX	SA	LCP	Before LCP (e.g. $SA[i],$ suffix $SA[i+1]$)
\$ \$	12	0	1
i \$	11	1	2
ippi \$	8	1	2
i ssippi \$	5	0	1
mississippi \$	1	0	1 2
P \$	10	1	1
ppi \$	9	0	3
Si PPI \$	7	2	2
SiSSippi \$	4	1	4
SiSiPPI \$	3	3	1
SiSiSSippi \$	2		
SiSiSiPPI \$	1		

SUBSTANTIAL SEARCH PROBLEM:

"Searching for a pattern P as a substring of T = **searching for pattern P as a PREFIX of a string in $SA(T)$.**"
(We have re-phrased our problem, after looking at $SA(T)[j]$.)

~~W~~ To solve this problem, we hence state 2 OBSERVATIONS:

1) SUFFIXES prefixed by P occur contiguous in $SA(T)$

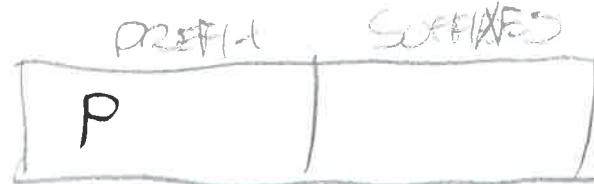
Ex: $P = "si"$

|
Sippi\$ } consecutive suffixes
Sissippi\$ } (with same prefix)
Ssippi\$

2) The lexicographic position of P in $SA(T)$ immediately precedes the block of suffixes prefixed by P .

Ex: $P = "si"$

|
si
Sippi\$
Sissippi\$
Ssippi\$



\Rightarrow These properties help in performing "PREFIX SEARCH" to achieve SUBSTRING SEARCH.

~~SUFFIXES~~

$LCP =$ Longest Common Prefix among two consecutive suffixes.

$SUFF[SAC[i]]$ and $SUFF[SAC[i+1]]$

BINARY SEARCH FOR
SUB-STRING-SEARCH PROBLEM:

We have reduced the problem of finding a SUBSTRING by finding a PREFIX \Rightarrow Hence, we can make use of Suffix String-Based Binary Search over SA/T.

SUBSTRING Search ($P, SA/T$)

$L=1, R=h$

while ($L \neq R$) do

$M = L \lfloor \frac{L+R}{2} \rfloor$

if ($strcmp(P, SUFFM, P) > 0$)

$L = M + 1$

// Now need to go right
of M

else:

$R = M$,

// need to go left of M

end if

end while

if ($strcmp(P, SUFFL, P) = 0$):

return L;

else

return -1;

end if



The only difference lies in WHERE we apply the APL. / here, we are applying it over a SUFFIX ARRAY \Rightarrow Go UP / Down along it).

BINARY SEARCH over SA \Rightarrow $O(\log n)$ = #comparisons

S consists of a string \Rightarrow $O(P)$ (worst case)
(for applying strCmp)

\Rightarrow **TOTAL TIME** = $O(P \cdot \log n)$ /
(SEARCH + COMPARISON)

~~We also need to count #occurrences~~

#occ: $O(P \cdot \log n)$
 P

strCmp \Rightarrow scanning of suffix requires
all $O(P)$

To count the #occurrences. Oocc
 $\Rightarrow O(Occ) = #occ.$

SPACE = ~~O(n.P.n.Log n + n.P.log P)~~

SPACE = $n \cdot (\underbrace{\log n + \log P}_{SP(T)})$ BITS
 T

APPLICATIONS OF SUBSTRING SEARCH:
Bio-informatics / ~~bioinfo~~ DNA sequencing

EXAMPLE SUBSTRIN SEARCH ALGORITHM
→ BINARY-SEARCH-LIKE APPROACH

We know that: $P = "ssi"$

SORTED
 $SUF(T)$:

1 \downarrow \$
2 i\$
3 iippi\$
4 iissippi\$
5 i sissippi\$
6 \xrightarrow{m} mississippi\$
7 P\$
8 PP\$
9 SIPP\$
10 SISI PPI\$
11 SISSI\$
12 \xrightarrow{r} SSISSI\$

SORTED
 $SUF(T)$:



SORTED
 $SUF(T)$:

1 → pi\$
2 PPI\$
3 SIPP\$
4 SISI PPI\$
5 SISSI\$
6 \xrightarrow{m} SSISSI\$
7 SSISSI\$
8 SSISSI\$
9 SSISSI\$
10 SSISSI\$
11 SSISSI\$
12 \xrightarrow{r} SSISSI\$

$$\text{① } \text{strn cmp}(\cancel{\text{ssi}}, \text{mississippi}) = 1 \quad \cancel{\text{ssi}} \Rightarrow$$

$$\text{strn cmp}(\text{ssi}, \cancel{\text{mississippi}}) = \cancel{\text{mississippi}} 1$$

$$\text{strn cmp}(\text{ssi}, \text{SIPP}) = 1 \Rightarrow \text{set } l = m+1$$

$$\Rightarrow l = r$$

Resulting string:

$r, l \rightarrow \text{SIPP}i$$

Mr. RLCP)

≡

Left side computed for
(WTIME)

ALGORITHM IS PROVIDED (To obtain better time complexity) ~~to avoid the DABBLE~~
 ↳ The comparison between P and M does NOT need to start from their initial character!

We exploit lexicographic sorting of the suffixes and skip character comparisons that have already been carried out in previous iterations).

THREE HELP ARRAYS & "MEMORIZATION"

- $LCP[1, n-1]$
- $LLCP[1, n-1]$ } Defined for every triple L, M, R that may arise in the inner loop
- $RLCP[1, n-1]$

The values in these arrays are:

$$- LLCP[M] = LCP(SUFFIX_{SAC[L]}, SUFFIX_{SAC[M]})$$

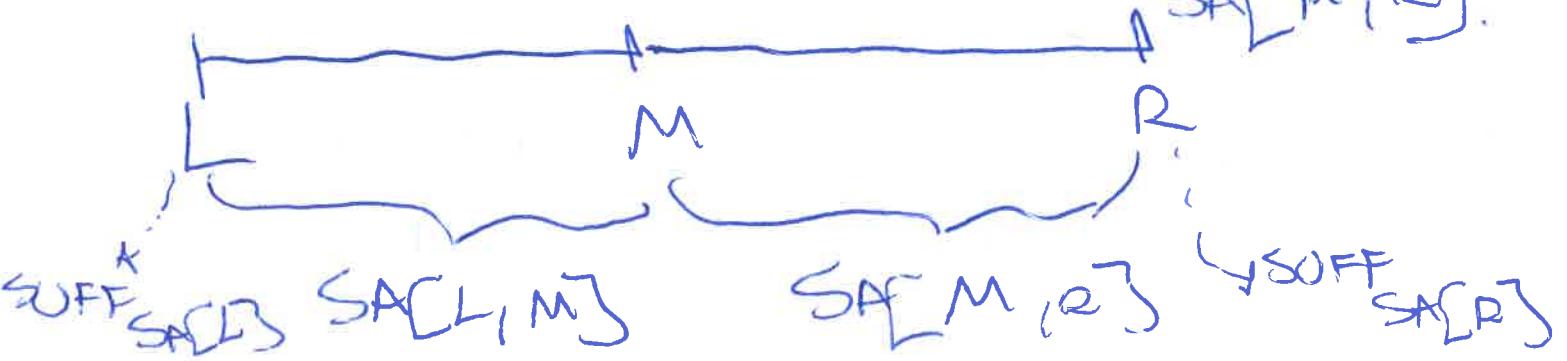
Prefix shared by the left-most suffix $SUFFIX_{SAC[L]}$ and right-most suffix $SUFFIX_{SAC[M]}$

$$- RLCP[N] = LCP(SUFFIX_{SAC[N]}, SUFFIX_{SAC[N]})$$

Prefix shared by middle suffix $SUFFIX_{SAC[N]}$ and right-most suffix $SUFFIX_{SAC[N]}$

(RLCP & LLCP can be built in $O(n)$ time via RMQ)

\Rightarrow At every iteration, we must decide whether P lies between $SAL[M]$ or $SAR[R]$.



\Rightarrow To decide where to go next, we need to perform a lexicographic comparison between $SAL[M]$ and $SAR[R]$.

(P) and SUFF $SAL[M]$

IDEA: Do not start from the first character every time! \Rightarrow would need $O(p)$ comparison!

\Rightarrow Take advantage of previously computed info.

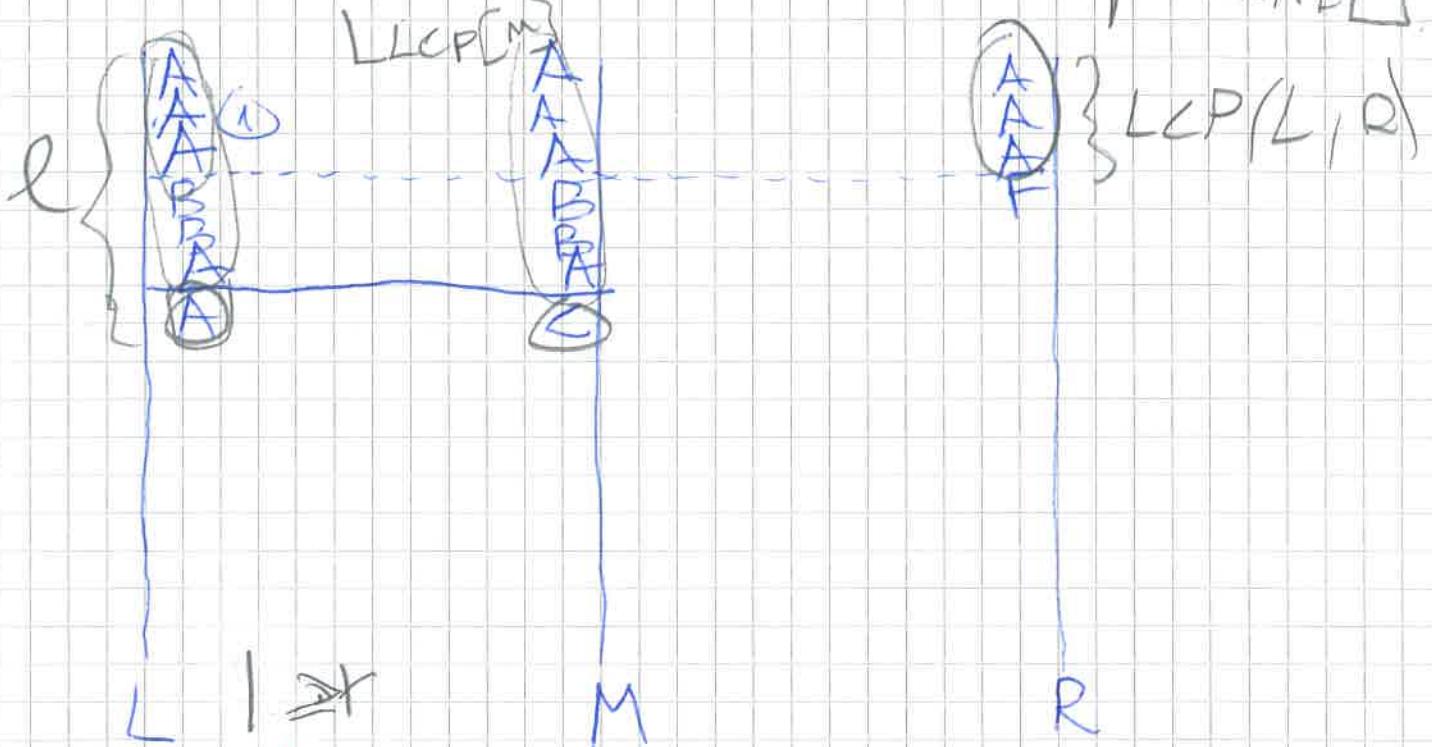
We know that:

P lies between $SUFF_{SAL[L]}$ and $SUFF_{SAR[R]}$

\Rightarrow $\text{LEP}(L, R)$ is the # characters that are shared by P with its suffixes (or all suffixes in the range) but share this # chars.

\Rightarrow We have learned how many characters P shares with L and R
 $(LCP[L, R])$

We can then draw 3 cases: $P = AAAABLM$



~~l = LCP(P, SUFF_{SAC[L]})~~

chars shared with left subarray

~~$l = LCP(P, SUFF_{SAC[R]})$~~

1) $l < LCP[M]$

\Rightarrow Then $P > SUFF_{SAC[M]}$, and we set

$m = l \Rightarrow$ Go RIGHT of M

By induction $P > SUFF_{SAC[L]}$, and their mismatching character lies at $l + 1$.

(e.g. $SUFF_{SAC[L]}$ shares more than l)
 character with $SUFF_{SAC[M]}$

\Rightarrow Mismatch between P and $SUFF_{SAC[m]}$
 Mismatch between "P and $SUFF_{SAC[L]}$

\Rightarrow Search continues on $SAC[M, R]$

(i.e.: we now need to go right of M)

\Rightarrow Binary search on the right of M .

2) If $L > LCP[M]$ \Rightarrow Binary search on the left of M .

P must be smaller than SUFF $SAC[M]$

$\Rightarrow m = LCP[M]$

\Rightarrow Search continues on $SAC[L, m]$

(i.e.: we now need to go left of M , with no more comparisons.)

3) If $L = LCP[M]$, then P shares L characters with SUFF $SAC[L]$ and SUFF $SAC[M]$

\Rightarrow Comparison between P and SUFF $SAC[M]$ can stop from their $(L+1)$ -th character.

\Rightarrow ADVANCE P !

\Rightarrow Every binary-search step either advances the comparison of P 's characters, or it does not compare any characters, but halves the range $[L, R]$ by performing a binary search over a new range. ($O(\log n)$)

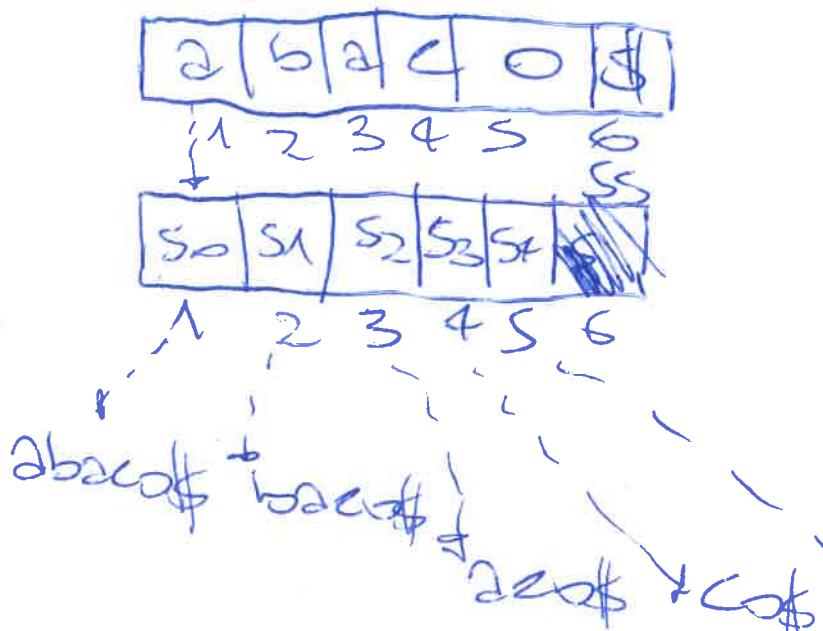
* TIME: $O(P + \log n)$ (count occurrences of a pattern P in file T)
for $STructure$

$O(P + \log n + \text{sec})$ (retrieve all the occurrences)

SPACE: $O(n)$

SUFFIX - TREY's Construction:

The suffix array is a sorted sequence of items \Rightarrow why not use a sorting algorithm that is applied over all the suffixes?



UNSORTED SUFFIXES		SA
a	bacott	1
b	acott	2
c	cott	3
d	t\$	4
e	ott	5

\Rightarrow We hence make use of a **comparison function + number server** where the comparison function is specialized to compute the lexicographic order between strings (e.g. arrays of characters) $\text{strcmp}(\text{SCFF}(X), \text{SCFF}(Y))$

\Rightarrow We do this in a pretty clever way, comparing character by character of different suffixes.

comparisons = $O(n \cdot \log n)$ \Rightarrow each string has n characters having $\log n$ comparisons per character

TIME = $O(n^2 \cdot \log n)$

I/Os = $O(n \cdot n \cdot \log n)$

SPACE = $O(n \cdot B)$ = $B \times \log n$ Space

```

SA_Construct(char* T, int n, char** SA)
{
    for (i=0; i < n; i++)
    {
        SA[i] = T + i; // new index add
    }
    QSORT(SA, n, sizeof(char*), suffix-cmp)
}

```

```

SUFFIX-CMP /char** p, char** q/
{
    return strcmp(*p/*q));
}

```

=> This procedure is not efficient
 at all though / not. #1/OS &
 time complexity)
 => STEW's algorithm is much
 more efficient!

CONSTRUCTION OF LCP ARRAY.

* The LCP is defined as:

$$LCP[i] = LCP(SUFFIX_{SA[i]}, SUFFIX_{SA[i-1]})$$

BRUTE-FORCE APPROACH:

$$SUFFIX_{SA[1]} = \underline{abaco}$$

$$SUFFIX_{SA[2]} = \underline{abaco} n$$

→ This approach takes $\Theta(n)$ for comparing 2 strings with one another.

→ For n suffixes, this approach takes $\Theta(n^2)$ time.

OPTIMAL LINEAR-TIME ALGORITHM:

TIME, SPACE: $O(n)$

→ We need to prove some other properties about Suffix Array. We need to avoid re-scanning the same text characters all the time to come to this result! (IS THIS POSSIBLE?)

~~By intuition from the figure that~~
~~Let's refer to a consecutive suffixes~~
~~SUFFIX ARRAY: suff_{i-1} and suff_i~~
~~at positions p and q in SA.~~
~~T = abcdefg~~

SOURCE SUFFIXES	SA	SA Positions	LCP
abdef	$\hat{S}-1$	P-1	\exists suffix?
<u>abchi</u>	$\hat{S}-1$	$P \rightarrow \text{suff}_1$?	
:			
bedef	\emptyset		2
bch	K	q-1	3
<u>bchi</u>	i	$q \rightarrow \text{suff}_2$	

Let's consider two consecutive suffixes (suff₁ and suff₂) at positions P and q.

$$T = abdefabchi$$

Assume: we know ~~intuitively~~ the value of LCP

$$\text{LCP}[P-1] = \text{LCP}(\text{SA}[P-1], \text{SA}[P])$$

\Rightarrow relevant for now test:

$$\text{LCP}[q-1] = \text{LCP}(\text{SA}[q-1], \text{SA}[q])$$

This can be computed without re-comparing the whole $\text{SA}[q-1], \text{SA}[q]$, but ~~the comparison can~~ part where ~~SA[P-1] < SA[P]~~ the comparison of $\text{SA}[P-1]$ and $\text{SA}[P]$.

We need 2 properties:

① $\forall x < y, it holds that:$

$$\text{LCP}(\text{SUFF}_{\text{SA}[y-1]}, \text{SUFF}_{\text{SA}[y]}) \geq \text{LCP}(\text{SUFF}_{\text{SA}[x-1]}, \text{SUFF}_{\text{SA}[x]})$$

(i.e.: any ~~preceding~~ preceding suffix to y will certainly share fewer characters than an adjacent suffix to x).
Proof: Prefixes are ordered!

JD MEANING: The longer (lexicographically greater) suffix from w₁ / The fewer ~~the~~ characters X shares with w₁

Ex. Consider suff_{w₁}¹ and suff_{w₁}².

base: suff_{w₁}¹ suff_{w₁}²

^{1b}
" ^{1b}
" " "

bcdef behi

^{1b}
" ^{1b}
" " "

$$\begin{aligned} \text{LCP}(p-i) &= \text{LCP}(\text{suff}^1, \text{suff}^2) \\ &\geq \text{LCP}(\text{suff}^1, \text{suff}^2) \\ &\quad \text{base} \quad \text{base} \\ &\quad \text{bcdef} \quad \text{behi} \end{aligned}$$

$$\Rightarrow 3 \geq 2$$

LCP_{w₁}

^{w₁ ends}
^{to w₁}

JD The following pair-wise suffixes will always be ~~the same~~ ^{w₁ ends} to w₁

LCP of the suffixes themselves.

② $\frac{w_1}{w_2}$

② EXPLAINED.

If SUFF_{j-1} and SUFF_{i-1} share some characters
(i.e. $\text{LCP}[P-1] \rightarrow \sigma$), then:

~~Blank~~ $\text{SUFF}_{j-1} < \text{SUFF}_{i-1}$ by the LEXICOGRAPHIC ORDER
 $\Rightarrow \text{SUFF}_{j-1} < \text{SUFF}_i$

And moreover, $\text{LCP}(\text{SUFF}_j, \text{SUFF}_i) = \text{LCP}[P-1]$
 \Rightarrow First shared character is dropped because
of the step ahead.

\Rightarrow The further (backwards) you go, the fewer the word characters will be.

(Proof: This property follows from the fact that suffixes must be ordered lexicographically.)

Follows from ①

$$\textcircled{2} \quad \text{by } \underline{\text{LCP}(\text{SUFF}_{\text{SA}[y-1]}, \text{SUFF}_{\text{SA}[y]}) > 0.}$$

Then: \Rightarrow

$$\begin{aligned} \text{LCP}(\text{SUFF}_{\text{SA}[y-1]+1}, \text{SUFF}_{\text{SA}[y]+1}) &\geq \\ \text{LCP}(\text{SUFF}_{\text{SA}[y-1]}, \text{SUFF}_{\text{SA}[y]}) - 1 &\quad \begin{array}{l} \text{PREVIOUS} \\ \text{SUFFIXES} \end{array} \\ &\quad \begin{array}{l} \text{NEXT} \\ \text{SUFFIXES} \\ \text{(decrement)} \end{array} \end{aligned}$$

MEANING: ①

$$\underline{abcdef} = \underline{\text{SUFF}_{\text{SA}[j-1]}}$$

and $\text{SUFF}_{\text{SA}[i-1]} = \underline{\text{abchi}}$

$$\underline{bcd\ell f} = \underline{\text{SUFF}_{\text{SA}[j]}}$$

\Rightarrow $\text{SUFF}_{\text{SA}[i]} = \underline{\text{bchi}}$

\Rightarrow Now take $\text{LCP}[P-1]$ of $\text{SUFF}_{\text{SA}[j-1]}$

• In the case that $\text{LCP}[P-1] > 0$: and $\text{SUFF}_{\text{SA}[i-1]}$

\Rightarrow Since $\text{SUFF}_{j-1} \ll \text{SUFF}_{i-1}$ (because $\text{LCP}[P-1] \leq \text{GRAPH}$)

\Rightarrow $\underline{\text{SUFF}_j \ll \text{SUFF}_i}$ and SORTING

$\text{LCP}(\text{SUFF}_j, \text{SUFF}_i) = \text{LCP}[P-1] - 1$ (DROP first char)

as we "STEP ahead" by 1 character
from \underline{P}_{i-1} to \underline{P}_i

\Rightarrow The # SUFFIX characters will be
 $(LCF[P-1] - 1)$
and the mismatching char. also remains.

Ex: In the case considered:

~~LCP[P-1]~~ & SUFF $SAC[i]$ and SUFF $SAC[j] = 3$
by
After grouping to SUFF $SAC[i-1]$ and SUFF $SAC[j-1]$
 $LCF[P-1]-1$ between $= \textcircled{2}$

\Rightarrow We have proved that the computation of $LCF[\alpha-1]$ can take advantage of what we have computed

at $LCF[P-1]$ (at the previous step) (O.E. D).

by
We AVOID RESCANNING!

TIME COMPLEXITY: $O(n)$

REQUIREMENTS:- Need the SA.
— Need all strings to be stored in INTERNAL MEMORY.

SUFFIX TREE: 10-16

It is a COMPACTED TRIE that stores all suffixes of a string T^{n+1} , where each suffix is represented by a unique path from the root to one of its leaves.

$$\# \text{LEAVES} = n$$

↳ leaves ~ representation (efficient)

↳ STRAIGHT PARENTS (STRAIGHT-POS, LEN) = $O(n)$

INTERNAL NODES = $\sum n$ \Rightarrow Its label signifies the # straight CROSSOVERS.

$$\begin{aligned} \# \text{EDGES} &= \text{TOTAL } \# \text{NODES} - 1 \\ &= \# \text{LEAVES} + \# \text{INTERNAL NODES} - 1 \quad \text{by} \end{aligned}$$

$\leq n + n - 1$ Each internal node has at least two outgoing edges!

$$\# \text{EDGES} \leq 2n - 1$$

↳ Edges are ordered alphabetically by their BLOCKING CHARACTERS.

(non-empty label starting the ~~suffixes~~ meaning the prefix stored by its children)

$$\text{SPACE} = O(n)$$

$$\Rightarrow \text{TOTAL SPACE (BLOCKS)} = O(n \cdot \log n)$$

PURPOSE of a SUFFIX TREE: It is used for efficient storage of suffixes and their efficient retrieval [quick substrace]

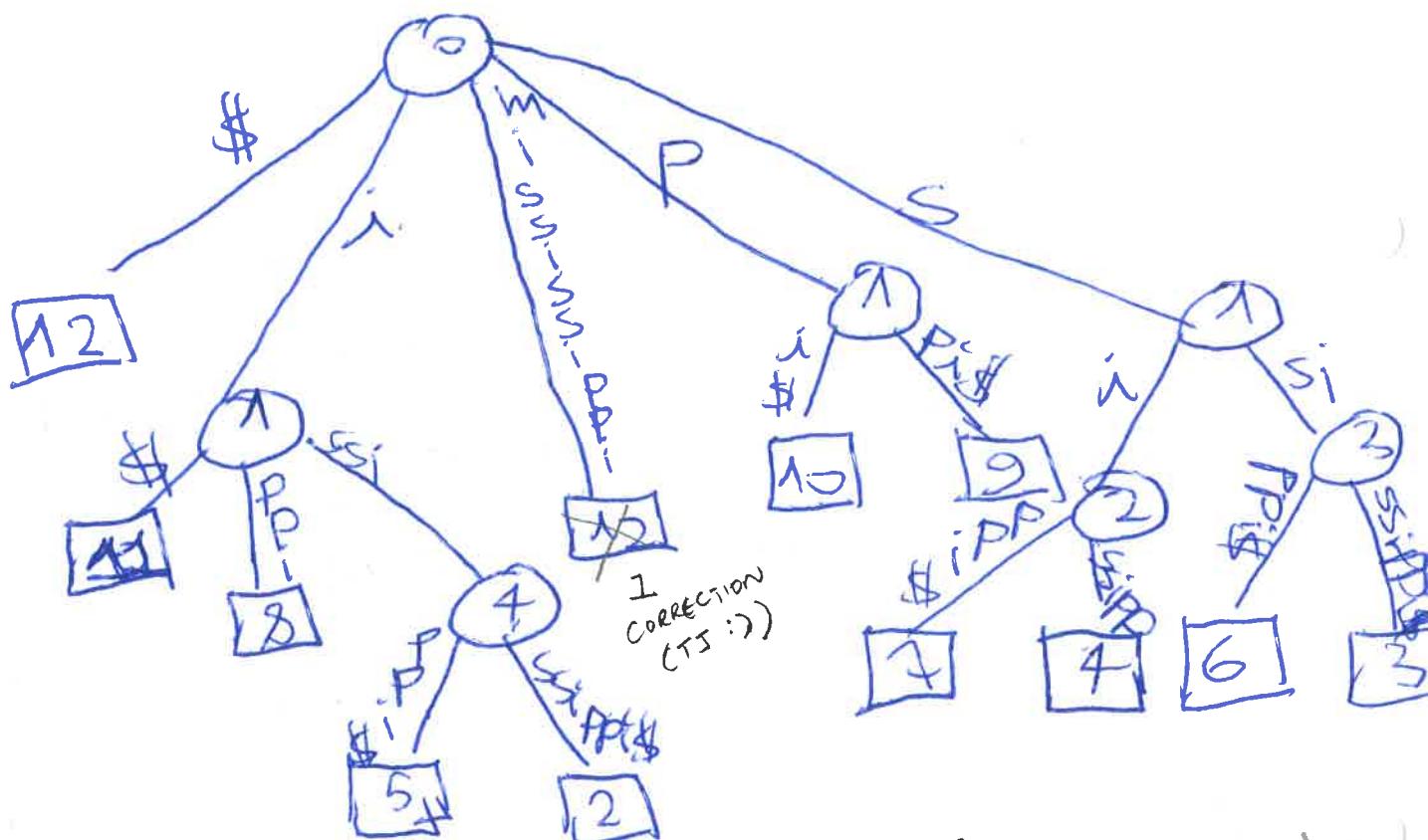
EXAMPLE SUFFIX TREE'S CONSTRUCTION.

Suffix Tree

SUF(T)	SA	LCP	LINEAR CONSTRUCTION
\$	12	0	1
i\$	11	1	2
ippi\$	8	1	1
issippi\$	5	4	4
ississippi\$	2	0	1
mississippi\$	1	0	1
pi\$	10	1	1
ppi\$	9	0	1
sippi\$	7	2	1
siissippi\$	4	1	1
ssippi\$	6	3	1
ssissippi\$	3	0	1

→ We get the preceding suffix of every single suffix!
by

Extrusion
LCP has
Same -
else
Common.



→ We use KMP's algorithm, we build the SA which return for every suffix its position in SA.
→ Set previous

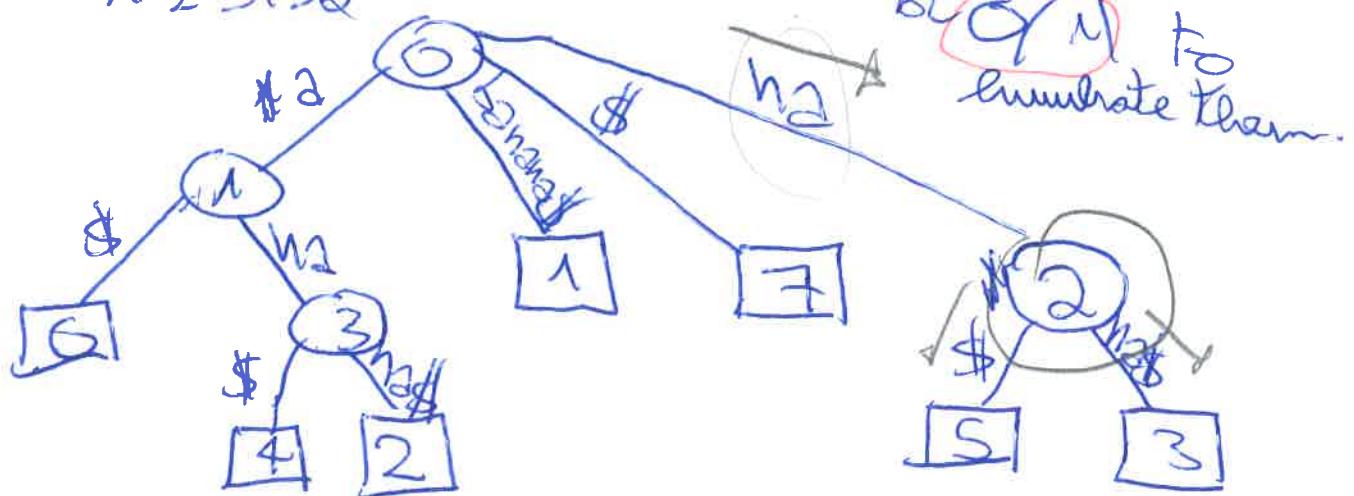
SUBSTRING-SEARCH PROBLEM.

We can use the ~~as~~-built SUFFIX TREE for pattern searching \Rightarrow How? We match characters along the way as the dog travels.

(only one downward path is possible!)

If the PATTERN is fully matched, then all nodes descending from the node found identify suffixes prefixes by $P \Rightarrow$ Office) to list

$T = b\text{anana\$}$



Ex: $P = n_2 \Rightarrow$ 2 occurrences Optimal

PATTERN SEARCHING cost: $O(p + n_{oc})$
Where $t_{oc} = \frac{\text{Time to branch out of a node}}{\text{during tree traversal}}$

t_{oc} can have different costs based on the implementation adopted.

① $T_0 = O(\log \theta)$ PLAIN ARRAY &
branching implements via BINARY SEARCH
Ex:

2	banana	#	n2	↓	2	b	\$		n		ab
---	--------	---	----	---	---	---	----	--	---	--	----

② $T_0 = O(1)$ as PERFECTION TABLE.

B	A	M	P	S
---	---	---	---	---

SPACE in ① | ② as OPTIMAL = $O(n)$

CONSTRUCT SUFFIX ARRAY & LCP FROM THE SUFFIX TREE:

We first define the:

- $\text{LCA} = \boxed{\text{Lowest Common Ancestor}}$ (of 2 leaves)

The Lowest Common Ancestor is the lowest node in the tree, which is a parent of multiple generations too & both LEAVES.

In the case of a SUFFIX TREE, we have that:

$$\boxed{\text{LCA}(\text{node}) \equiv \text{LCP}(\text{node})}$$

lowest common ancestor among
2 leaves

longest common prefix
among 2 leaves

LCP AND

- CONSTRUCT SUFFIX ARRAY From SUFFIX TREE.

The suffix array (SA) of text T can be constructed based on the SUFFIX-TREE by performing an IN-ORDER-VISIT of the SUFFIX-TREE.

Whenever a LEAF is encountered, its SUFFIX-INDEX is written to the SUFFIX ARRAY, along with the SUFFIX ITSELF.

TIME for SA

construction

$$= \boxed{O(n)}$$

SPACE

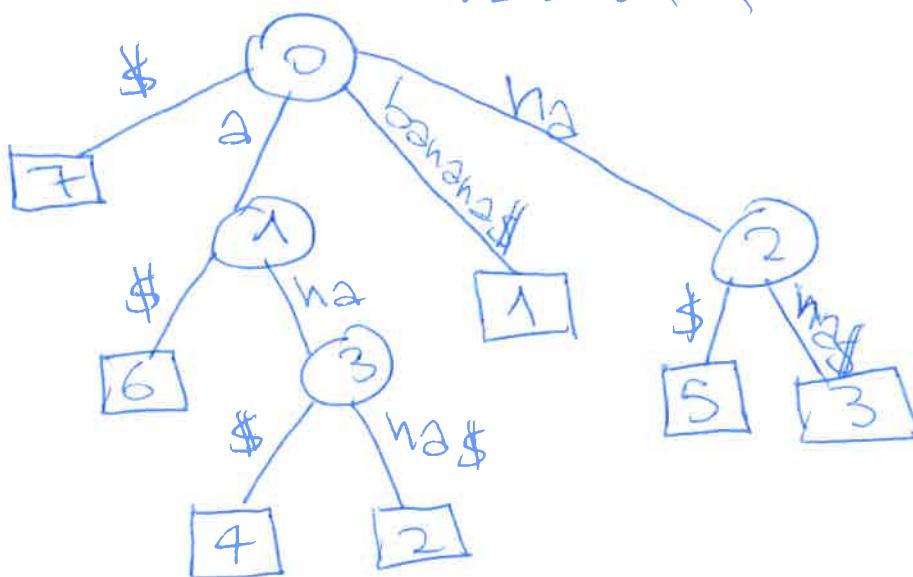
$\Rightarrow \text{SA} + \text{ST}$

$$= \boxed{O(n)}$$

~~EXAMPLE~~

ISLCP: During the In-order visit, we can write the value of an ~~already~~ encountered node in the LCP of the next visited node.

EXAMPLE, T = banana \$



→ Perform an IN-ORDER VISIT:

SEARCH SUF(i)	LEAVES	LCP	# CONS END	PREFIX
\$	\$	0	1	\$
a\$	7	1	2	a
an\$	6	3	3	ana
ana\$	4	0	1	
anana\$	2	0	1	
banana\$	5	2	1	ba
banana\$	3	0	1	

ALTERNATIVE APPROACH for LCP construction.
If we have a DATA STRUCTURE that allows us to perform LCA in $O(1)$ TIME, then we can compute the LCA of all pairs with nodes, i.e.

$$\text{LCA} \equiv \text{LCP}$$

Construct SUFFIX TREE from SUFFIX ARRAY: (& LCP)

INPUT: SUFFIX ARRAY

OUTPUT: SUFFIX TREE.

"DUMB" APPROACH: Take every single entry of the ~~SUFFIX ARRAY~~ & insert it them into the SUFFIX TREE $\Rightarrow O(n^2)$ time!

SMARTER APPROACH: Make use of the LCP computed so far, avoiding the need to compute it every single time for each pair.

We notice the following:

$$\text{LCP}[i] = \text{LCP}(\text{SA}[i], \text{SA}[i-1])$$

i.e.: The ~~LCP~~ ^{common suffix} of $\text{SA}[i]$ & $\text{SA}[i-1]$ is the LCA of SUFF $\text{SA}[i]$ and SUFF $\text{SA}[i-1]$.

$$\text{LCA}[i] = \text{LCA}(\text{SUFF } \text{SA}[i], \text{SUFF } \text{SA}[i-1])$$

IDEA: Compute the ~~consecutive~~ LCP/LCA between two entries of the SA (say $\text{SA}[i], \text{SA}[i-1]$)
IF $\text{LCP}[i]$ does not exist, then we invert the $\text{SA}[i]$ and $\text{SA}[i-1]$ underneath it.
(eventually percolating listing node)

IF $\text{LCP}[i]$ does NOT exist, then we invert ~~the SA~~ put $\text{LCP}[i]$ at the right place (percolating listing node)
& $\text{SA}[i]$ & $\text{SA}[i-1]$ underneath it. RIGHT PLACE (percolating listing node)

Ex $T = \text{banana\$}$

SORTED Suffix	SA	LCP
$\begin{array}{l} \$ \\ 2\$ \\ \text{ana\$} \\ \text{ananas\$} \\ \text{banana\$} \\ \text{hat\$} \\ \text{nana\$} \end{array}$	$\begin{array}{l} 7 \\ 6 \\ 4 \\ 2 \\ 5 \\ 3 \\ 0 \end{array}$	$\begin{array}{l} 0 \\ 1 \\ 3 \\ 0 \\ 2 \\ 0 \end{array}$

\Rightarrow Construct SUFFIX TREE based on SA.

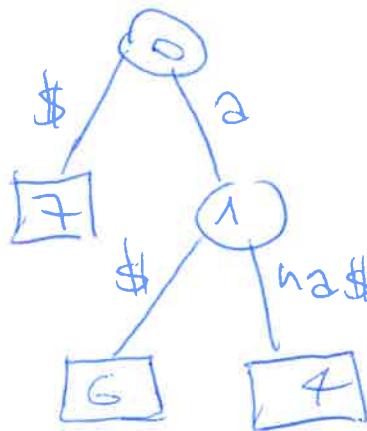
(1) $LCP = 0$

$$\begin{aligned} & SA[i] = 6 \\ & SA[i-1] = 7 \end{aligned}$$



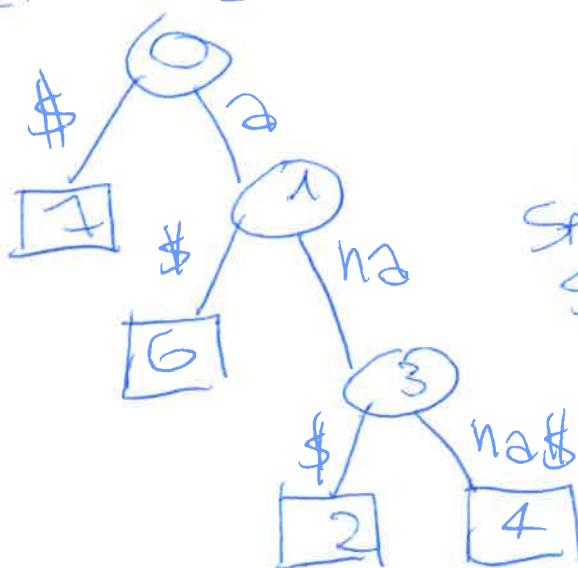
(2) $LCP = 1$

$$\begin{aligned} & SA[i] = 4 \\ & SA[i-1] = 6 \end{aligned}$$



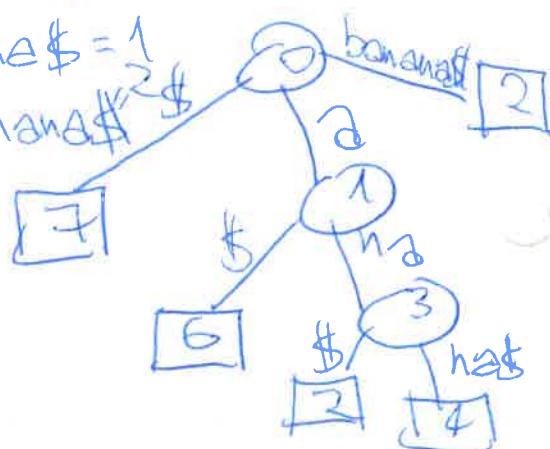
(3) $LCP = 3$

$$\begin{aligned} & SA[i] = \underline{\text{ananas\$}} = 4 \\ & SA[i-1] = \underline{\text{ana\$}} = 2 \end{aligned}$$



(4) $LCP = 0$

$$\begin{aligned} & SA[i] = \text{banana\$} = 1 \\ & SA[i-1] = \text{ananas\$} \end{aligned}$$



SEARCHING IN A SUFFIX TREE - PROPERTIES

① LCP between $SA[i]$ and $SA[j]$.
 $\min(LCP[i, j])$

WHY? Let's visualize the situation:



AAAAB B

AAAAC C

AAB

AA

AA.....

N.B.: No property forced actually!

→ Because of "TRANSITIVITY" & SORTING, we do have this ~~PROPERTY~~ property to hold.

Ex: LCP | SA | SUF(d)

LCP	SA	SUF(d)
5 → 6	12	\$
	11	i \$
	8	ippi\$
5 → 4	5	issippi\$

② \exists a SUBSTR of any length
that repeats \square times? ($l = 1$)

(Do we have a node from last
we have \square describing LEAVES?)
TDLNEx VISIT To find

③ \exists a SUBSTR of length ≥ 1
that repeats? (at least 2x)

Ex: ~~Want~~ Want a SUBSTR of length
at least 4.
MATCH



i.e.: do we have an ~~interval~~ interval
node ~~repeated~~ at least with 4?

? TDLNEx = LCP ≥ 4 ?

? Interval node!

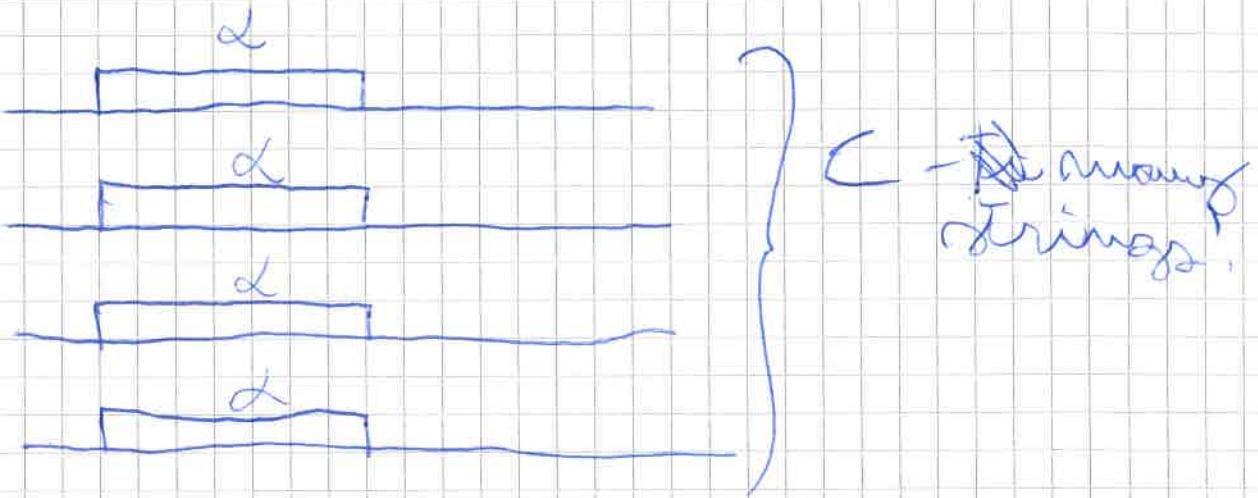
X a + ? } Strings share at least
a + ? } a characters!

Y a +

\Rightarrow Take all strings where LCP $\geq k$

④ \exists a SUBSTR of length ≥ 1
that repeats \square times?

\Rightarrow We have Θ strings starting with $\alpha = 1$.



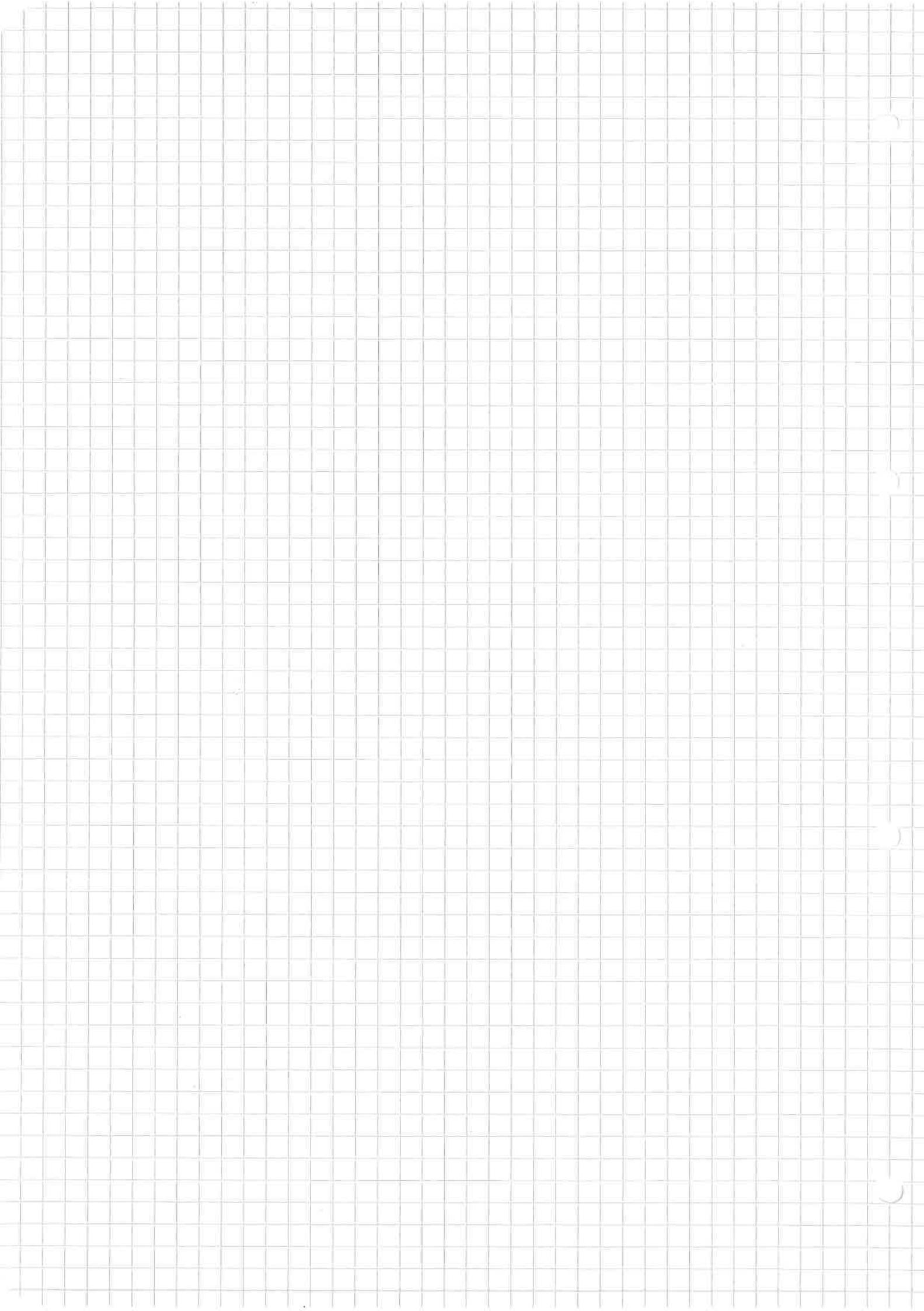
Storage for a Block of $C-1$ strings,
where:

$$\text{LCP} \geq L. \quad \text{(Last part } C-1\text{)}$$

Ex. APPLICATION: Intrusion detection system
~~System logging~~ based on system calls.

LCP	SA	SUF
0	12	\$
1	11	i\$
1	8	ippi\$
4	5	issippi\$
0	2	issippi\$

Ex: long 23
accru 22



CLASS EXERCISE

$T = \text{ananas} \$$

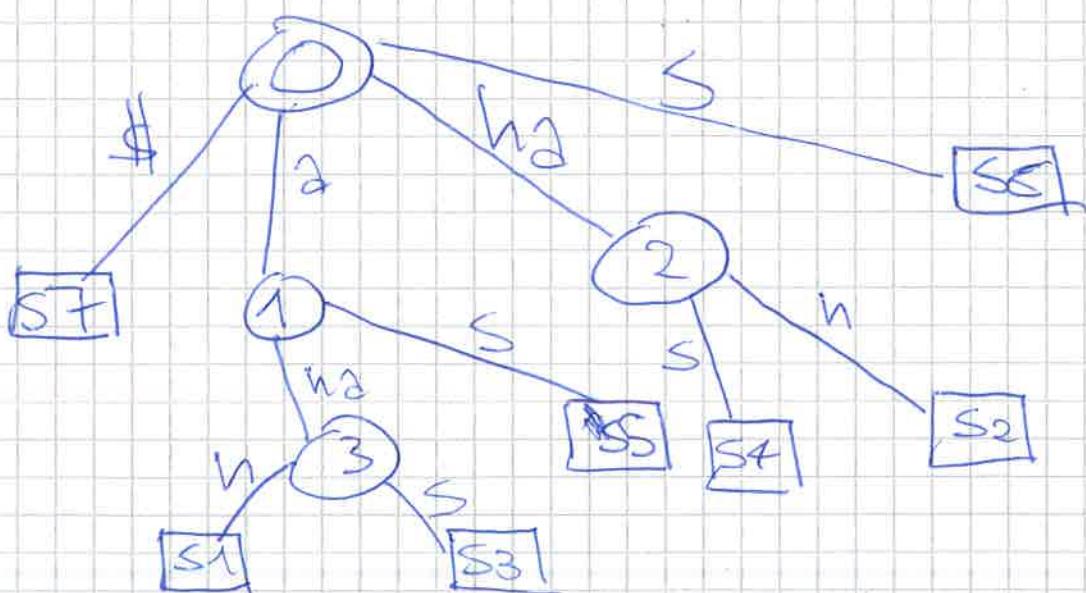
1 2 3 4 5 6 7

① Build a SUFFIX ARRAY.

SA	SUF(π)	LCP	# COMPS. (LINEAR)	# COMPS. (BUKE-FORCE)	LCP min
1	\$	0	1	1	0
2	ananas\$	3	7	4	2
3	nas\$	1	1	2	2
4	as\$	0	1	1	0
5	ana\$	2	1	3	1
6	as\$	0	1	1	1
7	s \$	0	1	1	0
NUMBER OF PREFIXES			16	14	

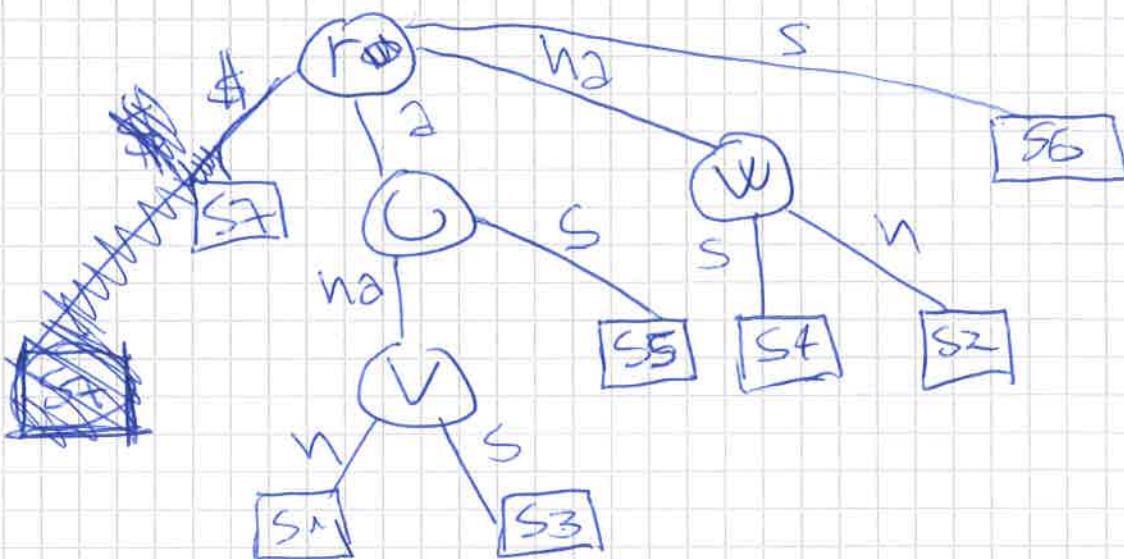
Store previously - found shared
PREFIXES (hence, the more
repetitive strings are, the
more comparisons we save)

② Build a SUFFIX TREE.



→ We now want to be able to solve RMQs - Range Minimum Queries over the SUFFIX TREE.

1 Build an EULER TOUR



2 Make an EULER TOUR over the tree:

Ex: ~~R~~ HUVWV3VWSURWVW2WV6
 Df: 0 1 0 1 2 3 2 3 2 1 2 1 0 1 2 1 2 1 0 1
 Slav: 0 + - + + + - - + - + - - - -)
 MIN: 3 2 1 ①

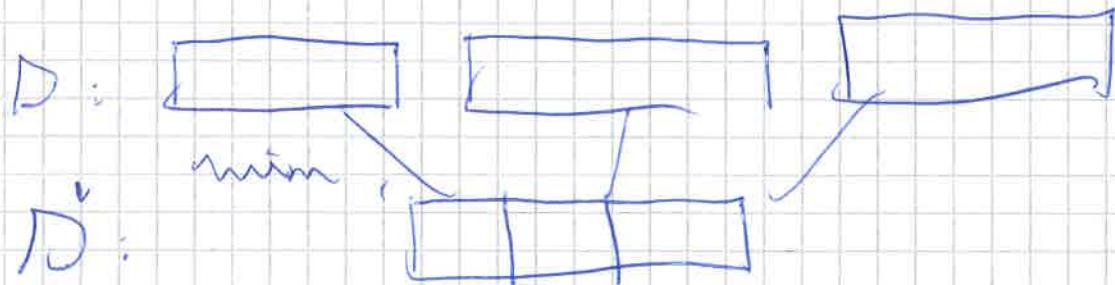
$$\text{BUCKET SIZE} = \log_2 \frac{n}{B} = 3$$

⇒ Gray. find LCA(1, 6) ⇒ solve RMQ over the depth.

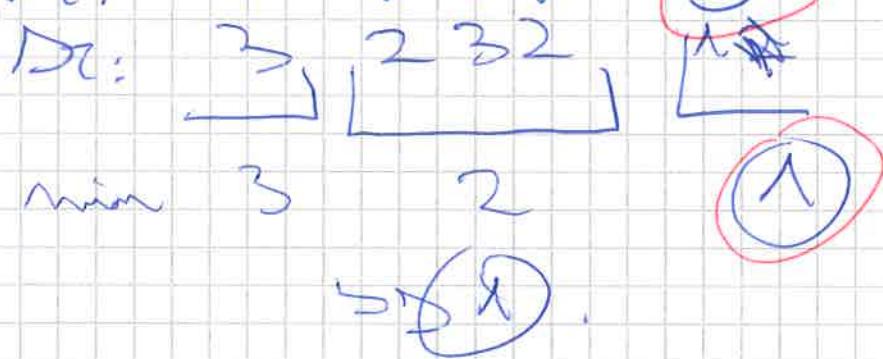
Build a sparse table, which contains the MIN. of ranges that are powers of 2.

$$\min_{i \in K} D^* [x_i, x_{i+2^K} - i]$$

NO NEED for a way inside a block.



Ex: \Rightarrow we compute it here in this way:



$$\text{LCA}(n_1, s) = 1$$

APPROXIMATE PATTERN MATCHING

PROBLEM STATEMENT:

"Find all substrings of a text $T[1, n]$ that match a pattern $P[1, p]$ with at most k errors"

Ex: $T = cabal[a]$

$P = bat[a]$

$\Rightarrow k=0 \Rightarrow$ No occurrence of P found in T .
 $k=1 \Rightarrow$ 1 occurrence of P found in T .
 k -Mismatch | (e.g. 1 error occurring).
 \Rightarrow Errors ~~Substitution~~

NAIVE SOLUTION: (BRUTE-FORCE)

try to match pattern P with every possible substring of T , taking care of the requirement that have at most k -mismatches with the pattern.

Ex: $T = cabal[a] \# MISMATCHES$ $P = bat[a]$

All substrings of T :
By length: $\{S_1 = cab, S_2 = abat, S_3 = bat[a]\}$

P

\Rightarrow This approach takes $O(n \cdot p)$ TIME.

SUBSTRINGS

N.B.: Rather inefficient approach, as we need to fully compare the pattern every time.
 \Rightarrow NEW APPROACH.

Do not use previously computed INFO.

RANGE - MINIMUM Query (RMQ)

→ Approach using the K-mismatch
Space in $O(n \cdot k)$ time [Provided we can
solve LCA between
P and a substring of T
in $O(N)$ TIME]

STEPS INVOLVED:

- ① Build a SUFFIX TREE in $O(n + p)$ TIME from ~~T~~^{DATA STRUCTURE} $\#P\#$
- ② Build LCA (lowest Common Ancestor) in $O(n + p)$ TIME, allowing to find the LCA (x, y) in $O(1)$ TIME.
- ③ Try to align pattern P with a SUBSTANCE of T [l, n], so that y_1 or y_2 -th substance coincide with y_1 character mismatch.

Ex.:
↑ MISMATCHES
↓ MATCHES
Interleave!!!

→ The search like to:

— Compute pattern P and test T

requires in $O(n)$ time is to get $O(n \cdot k)$
→ I need to be able to compute LCP
in $O(n)$ time.

→ This is the last step.

\Rightarrow In order to compare pattern $P[i], P[j]$ and a substring $T[i, j]$ [after being aligned as by Δ_T]

$T = \text{CGGATACGATCAGTAA}$
 $P = \text{CGGATACGATCAGTAA}$

\Rightarrow We would like to perform pattern & text equality in $O(n)$ TIME, to get $\mathcal{O}(n \cdot k)$ TIME.

Observation:

~~Keep focus in T with $\leq k$ mismatched~~

If $T[i:i+l] = P[j:j+l]$ is one of the matching substrings^{with k errors}, Then:

- l is the LCP between $T[i, \dots]$ and $P[j, \dots]$.

EXAMPLE:

~~LCP($T[1, 14], P[1, 7]$) = CCG~~

\Rightarrow THE matching substring respecting the OBS. no.

~~$T[4, 4+3] = P[4, 4+3] \Rightarrow l = 3$~~

- $\text{LCP}(T[5, 14], P[5, 7]) = AC$

~~$T[6, 6+1] = P[6, 6+1] \Rightarrow l = 1$~~

Q: How to compute:

$\text{LCP}(T[i:i+l], P[j:j+l])$ in $O(1)$ TIME

$\Rightarrow \text{LCP}(T[7, 14], P[7, 7])$

→ We know that SUFFIX TREES and SUFFIX ARRAYS contain some information about the LCP!

→ We ~~wants~~ are interested in SUFFIXES of P and T though!

→ We construct the SUFFIX ARRAY or SUFFIX TREE over the string:

$$X = T \# P \$$$

where # is a character not occurring anywhere else.

→ Through this newly constructed data structure, we are hence able to perform LCP computation among SUFFIXES.



(LCA)

LOWEST COMMON ANCESTOR & LCP:

From this Figure, we notice straightaway that there is a strong relation between the LCP problem over X 's suffixes and the LCA in the SUFFIX TREE. → The same can be found in the SUFFIX ARRAY.

Ex: Want to find ~~LCP($X[0, X], X[5, X]$)~~

$$\text{LCP}(X[0, X], X[5, X]) = \text{LCA}(X[0, X], X[5, X])$$

~~T = cabala #bara#~~

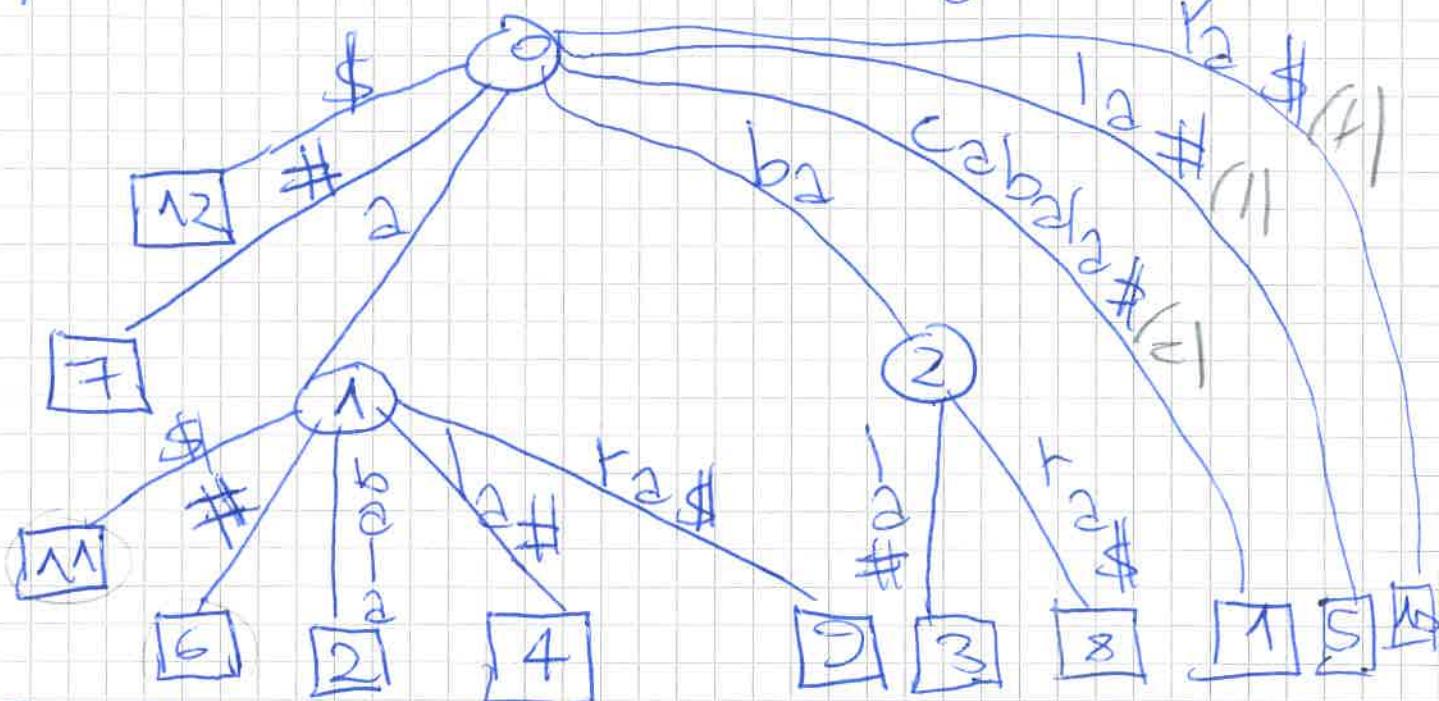
12 3 4 5 6 7 8 9 10 11 12

N.B.: ~~\$ - #~~

#SBS501.CS = 12

local term in
#PMSCTCCL
T-SB!

1) BUILD a SUFFIX TREE from T;



SORTED
SDF(T)

\$
#bara#
a\$
a#bara#
abala#bara#\$
ala#bara#\$
ara#\$
ba#bara#\$
bara#\$
cabala#bara##
la#bara#\$
ra#\$

SA:

12
7
11
6
2
4
9
3
8
1
5
10

LCP:

0
0
1
1
1
1
0
2
0
0
0
0

LCP(SUF[SAi, SAj])
SAi + j

(3) Provided we have a DATA STRUCTURE,
throughout we may find the LCA
in $O(N)$ time \Rightarrow We can:

~~Find the LCA among $\{S_1, S_2, \dots, S_k\}$~~
(& also the LCP's length!)

SUBSTRINGS

~~of S_1, S_2, \dots, S_k~~

~~To buy some toys, we must be
able to say that ~~exists~~ the
~~next~~ character from ~~small~~ will ~~be~~
at position ~~is~~.~~

EXAMPLE APPROACH

$T = C A B A L I X \#$
 $P = \underline{B P R A} \# \$$

\Rightarrow Look for the LCA of BARAH and
BALD # (in this manner
lets assume).

[8]

[3]

[6]

DS returns $\text{Q} = d$

~~BAKA #~~
~~BARA #~~

~~Same~~

2

[6]

Surely the character after the LCP
will differ!! ~~So find LCP~~

Find [] out []

~~Suppose compare A# and A# by
looking for their LCA~~

[6]

[We have a DATA STRUCTURE with that we can find the LCA in $O(1)$ time].

(1) APPROXIMATE: Solve in $O(nk)$ TIME via LCA computation over the SUFFIX TREE.

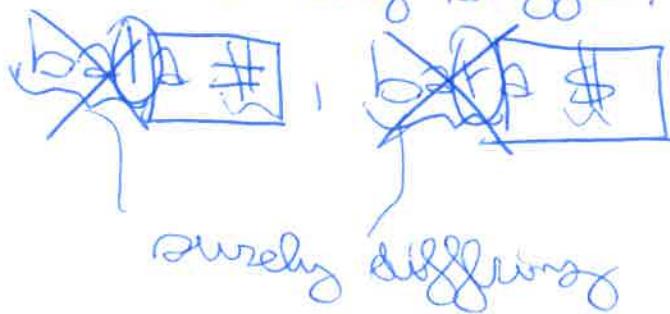
Ex: $T = cabab\#$; $K = 1$
 $P = \underline{bab}\underline{a\$}$

\Rightarrow We look for the LCA of ~~bab~~
and ~~bab\\$~~ in the MAGICAL LCA DATA
STRUCTURE.

(1)

$$\text{LCA}(\underline{bab}\#\underline{bab}\$) = 2 = \alpha.$$

\Rightarrow We know that the character following α will surely differ!



We are now left with performing one LCA query over $\underline{a\$}$ and $\underline{a\$}$

(2) $\text{LCA}(a\$, a\$) = 1$

\Rightarrow We have found that ~~T~~ and P share 3 characters, and the K=1 MISMATCH is satisfied.

~~NUMBER OF QUERIES performed = K+1~~

~~TIME COMPLEXITY~~ = $\mathcal{O}(n \cdot k)$

QUERIES over LCA = $\mathcal{O}(R)$

\Rightarrow The cost of these 2 comparisons is hence ①.

\Rightarrow We were looking for ① mismatch & we did 2 QUERIES
↳

In general, we make $k+1$ QUERIES over n characters at most.

$\Rightarrow O(n \cdot k)$ = TIME COMPLEXITY

#QUERIES over LCA = $O(k)$

② How to BUILD the LCA DATA STRUCTURE, to find LCA in $O(1)$.

[Among 2 different substrings].

RMQ: Range - Minimum Query PROBLEM.

Given an array $A[1:n]$ of elements drawn from an ordered universe, build a data structure RMQ that is able to compute efficiently the POSITION \min in $A[i:j]$

$O(n)$

$\forall i, j \in [1:n]$

POSITION \min

$\min \text{ in } A[i:j] \quad \forall i, j.$

Ex: A	1	2	3	4	5	6	7	8	9
	7	2	9	3	10	8	20	5	1

\downarrow

i

\downarrow

j

$\text{RMQ in } A[i:j] = 4$

BROKEN-FOOTER (TABLE-BASED) APPROACH:

(To get $O(1)$ TIME) in the SEARCH PHASE

- Store the index of a MIN. entry for each possible combination of i, j .
 (just pre-compute all answers)

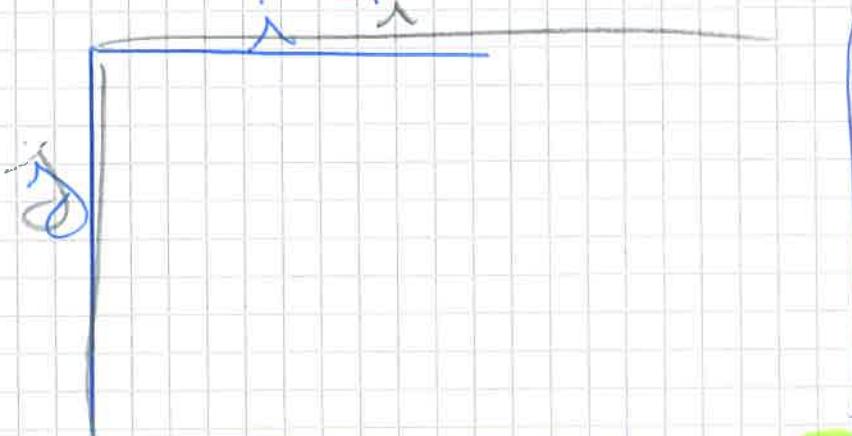


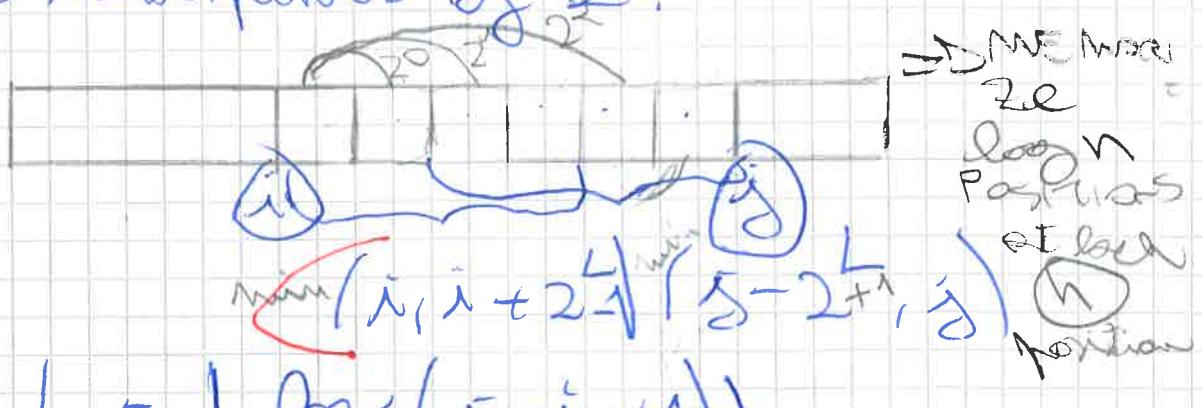
TABLE-build
time = $O(n^2)$

TABLE-SPACE
needed = $\Theta(n^2)$
Query time = $O(1)$

SPARSEIFICATION APPROACH. - $O(n \log n)$

Exploit ASSOCIATIVITY of the MIN. operations. We can hence try to decompose this DP ~~problem~~ & apply the MIN onto different pieces, then put them all together.

Observation: Any range (i, j) can be decomposed into two (possibly overlapping) ranges whose size is a power of 2.



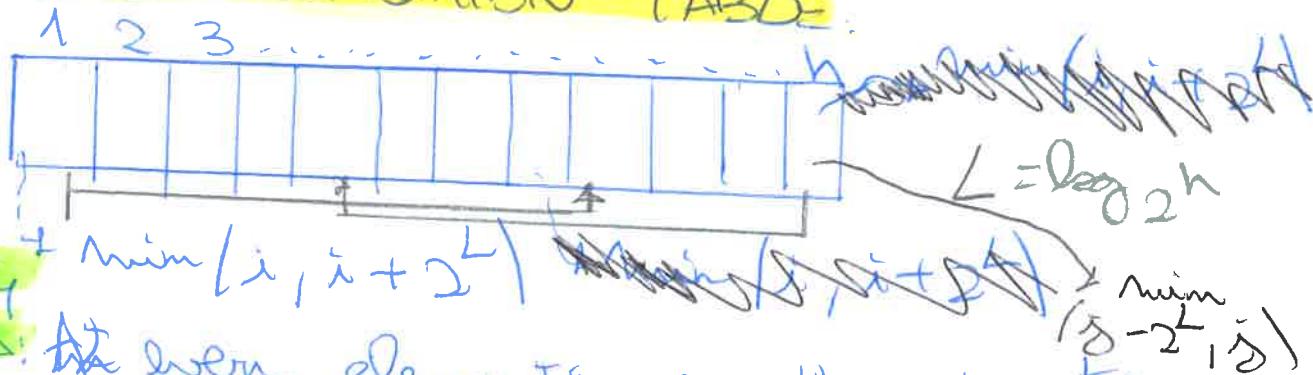
$$\text{Where } L = \lfloor \log_2(j-i+1) \rfloor$$

Idea: We "sparsefy" the quadratic-size table by storing only ranges whose size is a power of 2.

At position i , we hence store an answer to the query:

$$\boxed{\text{RMQ}(i, i+2^L)} \xrightarrow{\text{by storing}} \boxed{\min(i, i+2^L)} \text{ in}$$

the SPARSIFICATION TABLE.



I DED: At every element's position in the SPARSIFICATION TABLE, we pre-compute the result of the MIN over its logarithmic-sized ~~REPRESENTED~~ STORED ELEMENTS

We will then have pre-computed all values to answer any RMQ - QUERY.

$$\boxed{\text{RMQ}(i, j) = \min[\text{RMQ}(i, i+2^L), \text{RMQ}(j-2^L+1, j)]}$$

Where $L = \lfloor \log_2(j-i+1) \rfloor$

SPACE
(For SPARSIFICATION TABLE) = $\mathcal{O}(n \cdot \log n)$

QUERY TIME = $\mathcal{O}(1)$ (as we have pre-computed all possible results)
SPACE $\mathcal{O}(n \cdot \log n) \neq \text{OPTIMAL } \mathcal{O}(n)$ SPACE

OPTIMAL SOLUTION - TWO-FOLD REDUCTION:

In order to get SPACE = $\mathcal{O}(n)$ and

QUERY TIME = $\mathcal{O}(1)$

\Rightarrow We need to perform the following 2 Reduction's:

a) RMQ COMPUTATION over LCP ARRAY \rightarrow LCA COMPUTATION over CARTESIAN TREES.

b) LCA COMPUTATION over CARTESIAN TREES \rightarrow RMQ COMPUTATION over a Binary Array.

a) RMQ PROBLEM \Rightarrow LCA over CARTESIAN TREES

In this step, we construct the CARTESIAN TREE based on the entries of arrays $A[1, n]$ in a "RECURSIVE FASHION".

CARTESIAN TREE: Binary tree made of n nodes | where each node is labeled with:

$\boxed{\leftarrow \text{VALUE}, \text{position} \rightarrow}$

$A =$

7	2	9	3	10	8	20	5	1
1	2	3	4	5	6	7	8	9

 Different suffixes!

$\langle A[m], m \rangle$

$m = \text{Pos. of } \min(A[1, n])$
 $\forall A[m]$

LEFT SUBTREE

CART($A[1, m-1]$)

②

RIGHT SUBTREE

CART($A[m+1, n]$)
~~MIN~~

③ Built on the

right-hand side
of the min in
 $A[m]$

Built on the left-hand side
of the min in $A[m]$.

$\Rightarrow \boxed{\text{RMQ}_{\text{LCP}}(i, j) = \text{LCA}(i, j)}$ over the
CARTESIAN TREE.

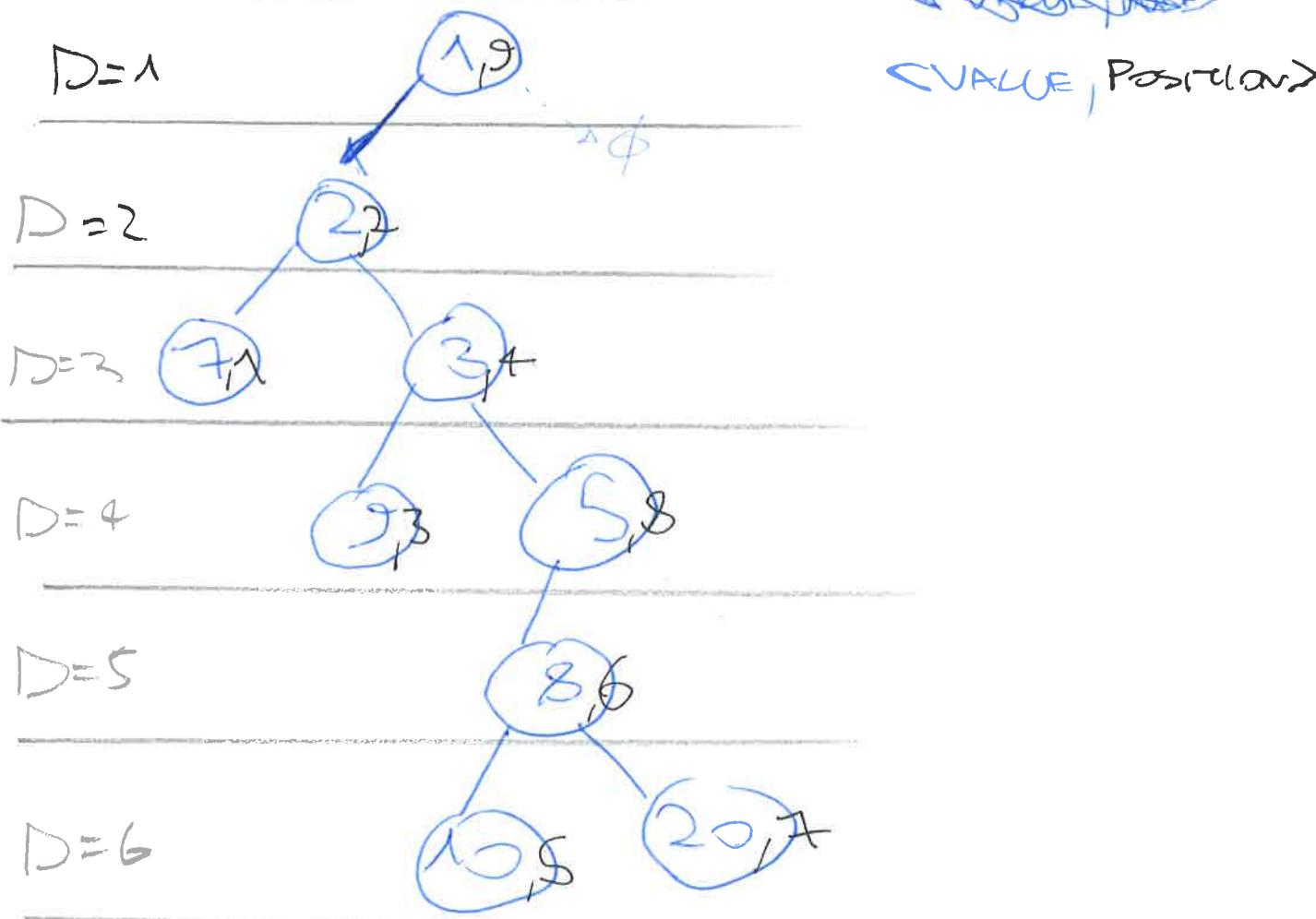
NB: (queried nodes i and j) may be internal
nodes as well

EXAMPLE - CARTESIAN TREE'S CONSTRUCTION

~~| * | * | * | 3 | 10 | 8 | 20 | \$ | *~~
 1 2 3 + 5 6 7 8 9

~~VALDE~~

⟨ VALUE, Position ⟩



Ex. queries:

one by the leftmost minimum

$$\text{RMA}(10, 20) = 80 = \text{LCA}(10, 20)$$

$$\text{RMA}(7, 3) = 22 = \text{LCA}(7, 3)$$

$$\text{RMA}(9, 8) = \langle 3, 4 \rangle$$

VALUE POSITION

\Rightarrow EULER TURE'S ORDINARY = $[2e + 1]$

b) LCA over (EXTREMA of BLOCKS only) / ADJACENT BLOCKS
 CARTESIAN TREE \Rightarrow RMQ PROBLEM over a BINARY ARRAY

\Rightarrow We build a BINARY ARRAY, which will allow us to (finally) solve the RMQ problem with $O(n)$ SPACE.

We do this according to the following steps:

- 1) Compute the EULER TOUR of the CARTESIAN TREE, writing down value when it's visited.
 (i.e. we perform a PRE-ORDER visit)
 A node can be visited multiple times too!
- 2) From the EULER tour E_T , build the DEPTH array D_T , which stores the depth of each visited node of the Euler tour.

$$\boxed{\text{LCA}(i, j) = \min_{i' < i} D[i', j]}$$

OVER THE CART. TREE

\Rightarrow Node of min. depth between i and j .

i' = pos. of first occurrence of i in the Euler-tour.

j = pos. of last occurrence of j in the Euler-tour.

i' = pos. of last occurrence of i in the Euler-tour.

~~This transformation allows us to re-write:~~

$$\boxed{\text{LCA}(i, j) = \text{RMQ}_D(i', j)}$$

SPACE for EULER tour + DEPTH array = $O(n)$

TRANSFORMATION TIME = $O(n)$, by doing first i 's last occurrence

- ~~3) Build Tree + "out" - "in"~~
- 3) We notice that the DEPTH array D_T has UNITARY VARIATIONS only (hence two consecutive nodes are connected by an edge).

→ PARTITION D_T into BLOCKS of size b .

where $b = \log_2 n$ $\Rightarrow \# \text{Blocks} = |D_T| = \frac{n}{b}$

For local blocks, we store the MIN. DEPTH together with its position.

$|M| = \frac{n}{b}$

$M[k]$ contains MIN. B_k of D_T .

(M_1, P_1)	(M_2, P_2)	(M_3, P_3)	(M_4, P_4)
--------------	--------------	--------------	--------------

$M_K = \text{MIN. of BLOCK } K$
 $P_K = \text{Pos. of MIN. in block } K$

4) Build on M the SPARSE-TABLE SOLUTION which solves RMQ queries in $O(1)$ TIME.

NB: M 's size is $\ll O(n)$.

M :

$\min(j-2^L+1, i)$

$\min(i, i+2^L-1)$

Where $L = \log_{\frac{1}{2}} M'$

SPACE For M and M' (Sparse Table)

$$O(|M| \cdot \log |M'|) = O\left(\frac{n}{b} \log \frac{n}{b}\right)$$

$$= O\left(\frac{n}{\log n} \log \frac{n}{\log n}\right) = O\left(\frac{n}{\log n} (\log n + \log \log n)\right)$$

$$b = \log n \Rightarrow O\left(\frac{n \cdot \log n}{\log n} + \frac{n \cdot \log \log n}{\log n}\right)$$

$$\Rightarrow O(n)$$

Ex. If $\text{RMA}(i, j)$ is performed over i, j belonging to different blocks. \rightarrow Doesn't work for

	B_1	B_2	B_3	B_4	B_5
E_T :	1 2 7 2	3 9 3 5	8 10 8 20	25 32 1	1 10 1
D_T :	1 2 3 2	3 4 3 4	5 6 5 6	5 4 3 2	1 1 2 5

$M:$	$\langle 1, 2 \rangle$	$\langle 3, 2 \rangle$	$\langle 5, 2 \rangle$	$\langle 2, 3 \rangle$	$\langle 1, 2 \rangle$
------	------------------------	------------------------	------------------------	------------------------	------------------------

$M =$	B_1	B_2	B_3	B_4	B_5
	$\langle 1, 2 \rangle$	$\langle 3, 2 \rangle$	$\langle 5, 2 \rangle$	$\langle 2, 3 \rangle$	$\langle 1, 2 \rangle$

Our query is: $\text{LCA}(9, 5) \rightarrow \text{LCA}(i, j)$ querying over blocks $B_2 \rightarrow B_5$

\Rightarrow Builds the SPARSEIFICATION TABLE.

B_2	B_5
3	3

$$L = \log(2^{i-j+1})$$

$$\cancel{\text{depth}} = 1$$

$$\min M\left(i, i+2^{\lfloor \frac{L}{2} \rfloor}\right)$$

$$\cancel{\min} = \min M[2:3]$$

$$\Rightarrow 3$$

$$\min M\left(j-2^{\lfloor \frac{L}{2} \rfloor}+1, j\right)$$

$$\min M[4-2^{1-1}, 4] = 2$$

⇒ The answer to

$$\text{LCM}(9, 10) = 30$$

$$\min(B_2, B_3) = 3$$

~~IF i, j ARE IN THE SAME BLOCK:~~

Then, the previously described approach wouldn't work.

∴ We need some other DATA STRUCTURES

To answer RMQs in $O(1)$ time over $i_0, j_0 \in B_i$.

- 1) Build the Euler tour ET (as before)
- 2) Build the Depth tour DT (as before)
- 3) Build the min-value M (as before)
- 4) Build the "++--" array ST , assigning a value based on the value of the preceding element.

$$ST[i] = \begin{cases} "+" & \text{if } DT[i] > DT[i-1] \\ "-" & \text{if } DT[i] < DT[i-1] \end{cases}$$

* NB: Not for the first element in the block!

- 5) Build The BINARY FINGERPRINTS $F(B_i)$ (as before).

$$\begin{aligned} "+" &\Rightarrow 1 \\ "-" &\Rightarrow 0 \end{aligned} \quad \boxed{F(B_i)} \quad \forall \text{ block}(i).$$

6) RETURN PAIRS:

$\langle F(B_i), \text{Pos. min in } F(B) \rangle$

7) Build a LOOKUP TABLE indexed by all possible query offsets i, j_0 within the block $B_K \Rightarrow$ to be able to answer queries in $O(1)$ TIME.

$T[\lambda_0, j_0 | F(B), \text{Pos. min in } F(B)]$

~~SPACE DEPENDS~~

Possible Configurations of FINGERPRINTS -
Space requires:

$$\begin{aligned} \# \text{ Possible Configurations} &= 2^{b-1} = O(2^b) \\ &= O(2^{\log n}) = O(2^{\log(n)^{\frac{1}{2}}}) = n^{\frac{1}{2}} = O(\sqrt{n}) \\ b &= \frac{\log n}{2} = \text{size of Fingerprint} \end{aligned}$$

TOTAL Space occupancy (of Lookup Table)

$$O(\sqrt{n} \cdot \log n \cdot \log n) = O(n)$$

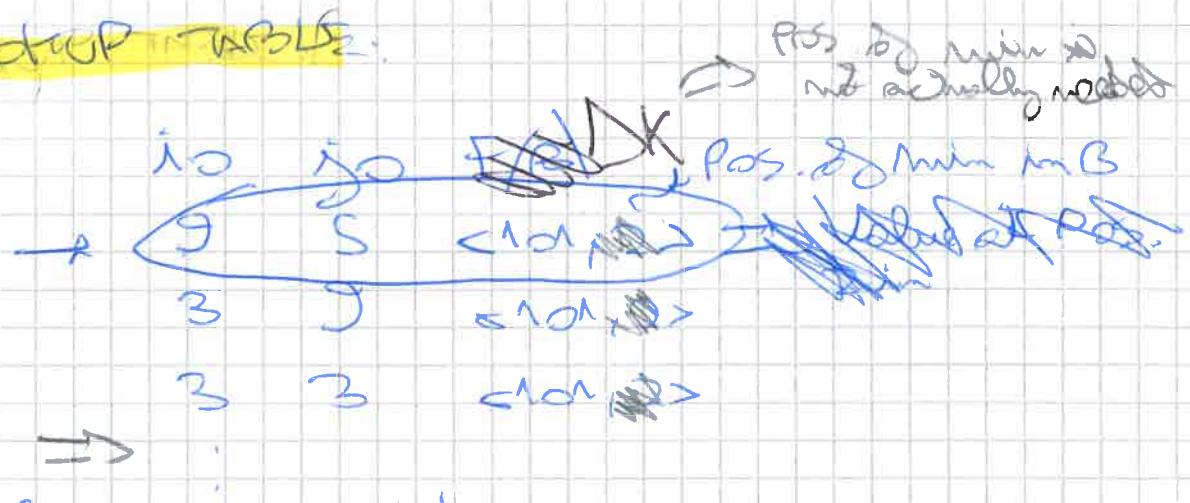
TIME to build the lookup table = $O(n)$

Final LookUp Time = $O(1)$

EXAMPLE of an RMQ over i_0, j_0 within the same block.

- 1) ET: $i_2 \overset{B_1}{\text{#}} j_2 \quad \boxed{3} \boxed{2} \boxed{3} \boxed{5} \quad 8 \overset{B_2}{\text{#}} 1 \overset{B_3}{\text{#}} 2 \overset{B_4}{\text{#}} 20 \quad 85321$
- 2) DT: $i_2 32 \quad 3434 \quad 5656 \quad 54321$
- 3) M: $\langle i_1, j_2 \rangle \quad \langle 3, j_2 \rangle \quad \langle 5, j_2 \rangle \quad \langle 2, 3 \rangle$
- 4) St: $i_1 + - \quad 3 + \cancel{-} + \quad 5 + - + \quad 5 - - -$
- 5) F(B): $\langle i_1 0 \rangle \quad \cancel{\langle 10 \rangle} \quad \langle 10 \rangle \quad \langle 000 \rangle$
- 6) PARS: $\langle i_1 0, 0 \rangle \quad \cancel{\langle 10, j_2 \rangle} \quad \langle i_1 0, j_2 \rangle \quad \langle 00, 3 \rangle$

Lookup Table:



\Rightarrow From $\langle F(B) \rangle$, we can hence answer the query
 $LCA(i_0, j_0) = B_2[2] = 3$.

ANSWER QUERY i_0, j_0 (in same Block)

1) Compute i_0, j_0 wrt. beginning of the block.
Keep below to.

2) ~~Rest~~ Compute $F(B)$ for that block.
 \Rightarrow Basis on $\langle i_0, j_0, F(B) \rangle$, the minimum hence be RETRIEVED.

CLASS
EXERCISE:

Let $T = abacabab$, $P = abt$.

a) How does the ϵ^1 -mismatch algorithm works on T and P , given that it is provided a DS that solves LCA in $O(n)$ time.

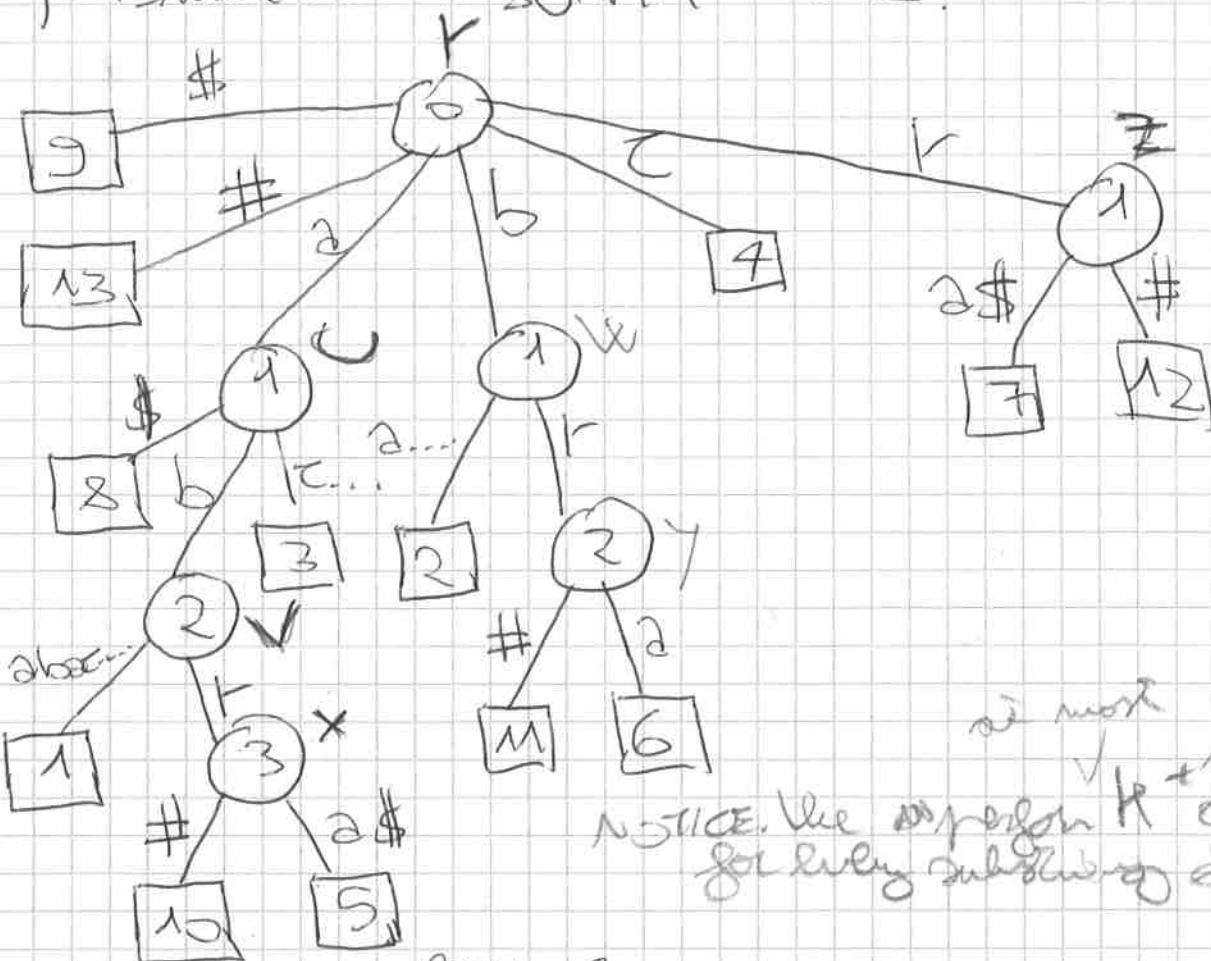
\Rightarrow Need to find 2 ϵ -mismatches (i.e. an LCA of T and P is enough.)

① Concatenate T and P .

P is enough!

$T \# P \# = \overbrace{ab@cabab}^{\text{LCA}} \# \overbrace{abt}^{\text{LCA}} \#$

2) Build the SUFFIX TREE.



\Rightarrow We compare P with the first suffix of T .

• $P \text{ COMP. } T[1: \dots] \Rightarrow \text{LCA}(1, 1d) = 2 \quad \checkmark \quad \epsilon^1 \text{ mismatch}$
 $P: (1, 12); T: (2, 12)$

• $P \text{ COMP. } T[2: \dots] \Rightarrow \text{LCA}(2, 1d) = \emptyset \quad \times \quad \{ \text{STOP}, 2 \text{ mismatch} \}$
 READING T: $\overbrace{abacabab}^{\text{suffixes}}$
 $\Rightarrow \text{LCA}(3, 12) = \emptyset \quad \times \quad \{ 2 \text{ mismatch} \}$

P comp. $T[3:\dots] \rightarrow \text{LCA}(3, 12) = 1 \Rightarrow \text{Next char. No a mismatch}$
 }
 Jump to next next char
 \downarrow
 $\downarrow \text{LCA}(5, 12) = 0 \Rightarrow \text{Swp. 2 mismatches.}$

P comp. $T[4:\dots] \rightarrow \text{LCA}(4, 12) = 0$
 $\downarrow \text{LCA}(5, 11) = 0 \Rightarrow \text{Swp. 2 mismatches.}$

P comp. $T[5:\dots] \rightarrow \text{LCA}(5, 12) = 3 \checkmark$
 (No mismatch in
 $T: (\alpha, \beta) P: (\#, \beta)$)

b) Provide a DATA STRUCTURE that allows for LCA in $O(1)$ time

i) Perform FOLDS tour.

~~Eg: #RUBUV1VX10x5xV03UW2~~
 DT: black size = 4
~~B1 B2 | B3 B4 | B5 B6 | B7 B8 B9~~
 Eg: RUBUV1VX10x5xV03UW2
 DT: 012123434321210121232321010121210
 M: 0 2 3 1 0 2 0 0 0

\Rightarrow We need to solve $\text{LCA}(10, 6)$

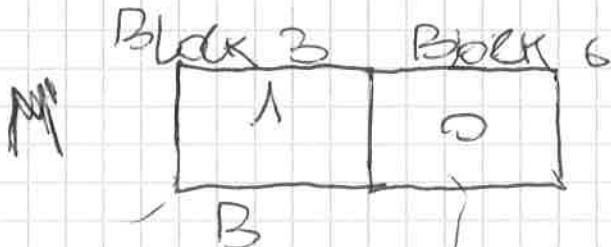
\Rightarrow Only check the SPARSIFICATION TABLE

Blocks	3	4	5	6	INVERSE VISIT
BR	B8	B5	B6	B3	
M =	3	1	0	2	

To answer the query given $\rightarrow \text{LCA}(10, 6)$

- ① It is hence sufficient to apply
SPECIFICATION to ^{the} blocks where ② and
③ lie in.

$$L = \log_2 (j-i+1)$$



$$\log_2 / 3$$

$$\min(j, j+2^L-1) = \min(3, 3+2^1-1) = 1$$

$$\min(j-2^L+1; j) = \min(6-2^1+1; 0) \\ x = 0$$

$$\Rightarrow \text{LCA}(10, 6) = \min(1, 0) = 0!$$

WORKS ONLY WITH
MULTIPLES OF BLOCKS

$$L = \lfloor \log_2 (j-i+1) \rfloor$$

$$L =$$

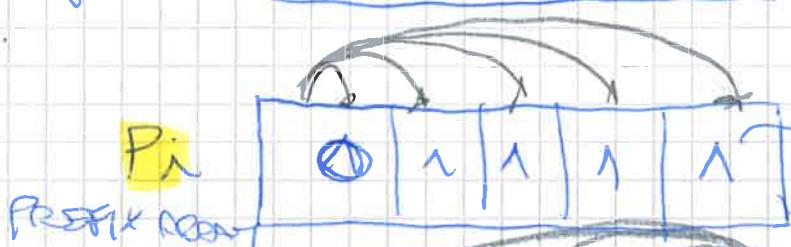
SPARSE TABLE - PREFIX / SUFFIX Array

This strategy allows us to answer queries among positions that are not multiples of blocks in O(1) time. It's the starting and ending position of a block, like in the following successive approach.

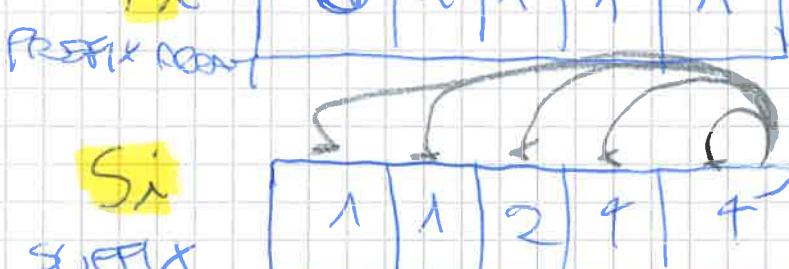
- 1) Build the FULL TREE (E-1) over the CARTESIAN TREE, as in 1)
- 2) Build the DEPTH TREE (as in 2)
- 3) Build the MIN-TABLE (M) using ranges in-between blocks
- 4) Build the CROSS-TABLE (M') for the ranges in-between blocks i and j
- 5) Build the PREFIX ARRAY & the SUFFIX ARRAY for the blocks that contain i and j

D_i	1	2	3	+	
	5	3	4	9	7

[ABSOLUTE POSITIONS]



Pos. of MINIMUM in $D_i[0 \dots i]$



[ABSOLUTE POSITIONS]

Pos. of MINIMUM in $D_i[i \dots n]$

$D_i[i \dots n]$

- 6) Take MINIMUM of the PREFIX ARRAYS (left part) SPARSE TABLE (middle part, OPTIMAL), SUFFIX ARRAY

EXAMPLE. i

E_T: 1 2 7 5
D_T: 1 2 3 2
M: 5 1 1 2
 $S_1 \begin{bmatrix} 2 & 2 \\ 3 & 3 \end{bmatrix}$

B₁: 39 35
B₂: 8 10 8 20
B₃: 5 6 5 6
B₄: 8 5 3 2
 $M' \begin{bmatrix} 3 & 0 & 3 & 0 \end{bmatrix}$

$\begin{array}{c} 3 \\ 8 5 3 2 \\ S 4 3 2 \\ \leftarrow 2 3 \right\rangle \\ S 4 3 \\ 9 9 1 \end{array}$
 $P_4 \begin{bmatrix} 0 & 1 & 2 \end{bmatrix}$

~~$S_1 \begin{bmatrix} 2 & 2 \\ 3 & 3 \end{bmatrix}$~~
 ~~$S_1 \begin{bmatrix} 1 & 1 \\ 3 & 2 \end{bmatrix}$~~

~~$S 4 3 2$~~
 ~~$P_4 \begin{bmatrix} 0 & 1 & 2 \end{bmatrix}$~~

$$\min(S_1, M', P_4) = 2$$

① INTEGER CODING.

PROBLEM: $S = \{s_1, s_2, \dots, s_n\}$ be a sequence of integers (positive), $s_i \leq N$, possibly repeat.

GOAL: Represent integers of S as binary sequences, what are SELF-DETERMINING and occupy few bits.

If integers are NOT POSITIVE
↓

We transform -POSITIVE INTEGERS \rightarrow EVEN Positive INTEGERS

-NEGATIVE INTEGERS \rightarrow ODD Positive INTEGERS

$$T \subseteq \mathbb{Z}$$

$$\boxed{x \in T \geq 0 \rightarrow 2x \text{ (EVEN)}}$$

$$\boxed{x \in T < 0 \rightarrow -2x + 1 \text{ (ODD)}}$$

N.B.: Some numbers can encode (6), others cannot.

FIXED-LENGTH REPRESENTATION

(i.e. 4/8 bytes to represent an integer)

Take $M = \max_{s \in S} s$ (i.e. max. integer in the set S)

$$\forall x \in S < M \Rightarrow \# \text{BITS} = 1 + \lfloor \log_2 M \rfloor$$

SIMPLE BINARY ENCODING! (Ex: 3 → 11)

GOOD WHEN?

BAD WHEN?

When integers not zero
sparse (containing many 0s)

Integers are ~~not~~ very sparse (too many 0s!!)
dense (M has 8 bits!!)

- BAD.

Not self-delimiting code if concatenating multiple integers' encoding [we wouldn't know where an integer starts and the next one ends at].

↳ NO UNAMBIGUOUS DECODING.

Ex: $S = \{1, 2, 3\}$

$C = \{110111\} \Rightarrow$ How do we decode this?
(DUNNO!)

VARIABLE-LENGTH REPRESENTATION

(all the following numbers)

We make life easier & succinct many of coding and decoding, which requires a variable #BITS for space & for UNAMBIGUOUS to be decoded.

WANT: OPTIMAL ENCODING.

We have an OPTIMAL ENCODING of integer X if

$$\sum_x P(x) |\text{ENC}(x)| \geq \sum_x P(x) \log_2 \frac{1}{P(x)}$$

Avg. #BITS used from encoder | Encoder (SHANNON)
↓ Lower-bound | Coding Theory
→ P. &
asym
x curve

Whatever the encoding, we can NEVER do better than this!
(No Proof!)

Provides an UPPER BOUND to the #BITS used to encode integer X .

SHANNON'S CODING THEORY.

The ideal code length $L(x)$ for a symbol x is: $\log_2 \frac{1}{P(x)}$ (Entropy) $P(x) = P\{x \text{ occurs}\}$

~~Upper Bound of H?~~

$$0 \leq H \leq ?$$

\downarrow (Upper Bound)

$$H = \sum_{x \in S} P\{x\} \cdot \log_2 \frac{1}{P(x)} \quad \text{where } x \text{ is a symbol in } S$$

$$\sum \sum_{x \in S} \frac{1}{|S|} \cdot \log_2 |S| = \log_2 |S|$$

Upper Bound
= Entropy

$$\Rightarrow 0 \leq H \leq \log_2 |S|$$

Ex: If you have 16 digits ($S=16$), the entropy will go from 0 to 4.

Entropy grows (slowly) with the # of bits. In sets,

Entropy's meaning: Avg. # Bits to encode a string (based on source) emitting x .

$\Rightarrow H \rightarrow 0$. Highly compressed,

} values are concentrated at some values (some NOT random)

$\Rightarrow \log_2 |S|$, sparse values.
} values are NOT concentrated

Random source.

Optimality Condition | Equation.

In order for an encoding scheme to be optimal, we must have:

$$|\text{ENC}(x)| = \log_2 \frac{1}{P(x)} = \underbrace{\text{#BITS required}}_{\text{for encoding } x} \quad \text{for an integer } x \text{ occurring in the distribution considered}$$

Avg. #BITS required

Ex:

$$S = \{1, 2, 3\}$$

P: $(0.5, 0.3, 0.2) = P\{x\}$

BAD Encoding:

$$\begin{aligned} 1 &\rightarrow 000 \\ 2 &\rightarrow 10 \\ 3 &\rightarrow 11 \end{aligned}$$

Lots of BITS!

$$\begin{aligned} |\text{ENC}(x)| &= 0.5 \times 3 + 0.3 \times 2 + 0.2 \times 2 \\ &= 1.5 + 0.6 + 0.4 \\ &= 2.5 \text{ bits} \end{aligned}$$

GOOD Encoding:

$$\begin{aligned} 1 &\rightarrow 0 \\ 2 &\rightarrow 10 \\ 3 &\rightarrow 11 \end{aligned}$$

$$\begin{aligned} |\text{ENC}(x)| &= 0.5 \times 1 + 0.3 \times 2 + 0.2 \times 2 \\ &= 0.5 + 0.6 + 0.4 = 1.5 \text{ bits} \end{aligned}$$

\rightarrow We hence notice that, as we perform a weighted graph based antile lexicographical length and the P. & an integer occurring DESIRED PROBABLY SYMBOL \rightarrow Long encoding PROBABLE SYMBOL \rightarrow Short encoding

OPTIMAL DISTRIBUTION for a FIXED-LENGTH ENCODER:

(i.e. for which distribution is ~~with~~ a fixed length encoder OPTIMAL?)

~~$$\forall x, \text{Length}(\text{Enc}(x)) = \log_2 \frac{1}{P\{x\}}$$~~

P(x) needs to encode $\log_2 \frac{1}{P\{x\}}$

$$|\text{Enc}(x)| = \log_2 \frac{1}{P\{x\}}$$

(NOT considering +1)

$$\log_2 M = \log_2 \frac{1}{P\{x\}}$$

$$P\{x\} = \frac{1}{M}$$

$\forall x \in M$

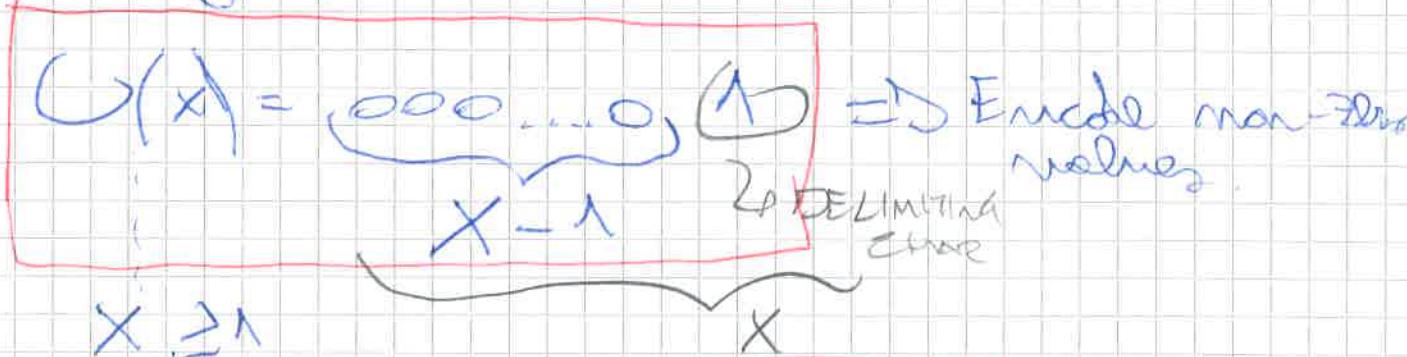
2 MAX. value
of size 5

\rightarrow FIXED-LENGTH encoder is optimal if the DISTRIBUTION is UNIFORM!

UNARY CODE: $\{x\}$

IDEA:

$U(x)$ for an integer $x \geq 1$ is given by a sequence of $X-1$ bits set to 0, unless there is a (DELIMITING) bit set to 1.



SPACE of $U(x) = \Theta(x)$ BITS!

(Exponentially more bits than BINARY CODE)
 ↗ FIXED-length representation

↗ Space inefficient for large x .

EXAMPLE: Ex: $U(5) = \underbrace{0000}_{4 \text{ '0s'}} 1$

OPTIMAL

DISTRIBUTION

OPTIMAL \Rightarrow UNIFORM ENCODING.

$$|\text{ENC}(x)| = \log_2 \frac{1}{P(x)}$$

As many bits as x value itself

$$\textcircled{X} = \log_2 \frac{1}{2 P(x)}$$

$$2^X = \cancel{\log_2 \frac{1}{P(x)}}$$

$$2^X = \frac{1}{P(x)}$$

$$\Rightarrow P(x) = \frac{1}{2^X}$$

18

EXPONENTIAL
DISTRIBUTION
OPTIMAL

ELIAS CODES:

δ and γ ENCODING are UNIVERSAL ENCODES, as their code length is $O(\log x)$ for any integer x , and prefix-free.

Also, they are UNAMBIGUOUSLY-DECODABLE, by inverting the Coding Procedure.

δ - (GAMMA) CODE:

Given an integer (x) as input, where $B(x)$ is its BINARY REPRESENTATION and $|B(x)|$ is the # bits used to represent $B(x)$.

$$L = |B(x)|$$

$$\delta(x) = \underbrace{000\ldots0}_{L-1 \text{ } 0s} \underbrace{1}_{\textcircled{1}} \cdot B(x) \underbrace{\overbrace{\quad\quad\quad\quad\quad\quad\quad}^{\textcircled{2}}}_{L}$$

$\delta(x)$ is formed by:

- 1) The unary representation of $|B(x)|$.
(i.e. $|B(x)| - 1$ zeros + 1 bit 1)
- 2) $B(x)$, the binary representation of x .

SPACE: $L-1 + L = 2L-1 = O(\log x - 1)$

EXAMPLE

$$0/0 \Rightarrow B(d) = 1001$$

$$0/0 = \underbrace{000}_{B/d-1} \underbrace{1001}_{\text{1 zero}}$$

OPTIMAL DISTRIBUTION for γ -ENCODING

$$I_{ENC}(x) = \log_2 \frac{1}{P\{x\}}$$

$$2 \cdot \log x - 1 = \log_2 \frac{1}{P\{x\}}$$

$$2 \cdot \log_2 x - \log_2 2 = \log_2 \frac{1}{P\{x\}}$$

$$\log_2 x^2 - \log_2 2 = \log_2 \frac{1}{P\{x\}}$$

~~$$\log_2 \left(\frac{x^2}{2} \right) = \log_2 \frac{1}{P\{x\}}$$~~

$$P\{x\} = \frac{2}{x^2}$$

DISADVANTAGES:

Need to perform Bit-By-Bit decoding & encoding "Slow"

↳ Actually used when small values are involved!

S-Encoding (DELTA)

It applies Δ -code instead of unary encoding.

$S(x)$ also consists of 2 parts:

- 1) Encode $\delta(|B(x)|)$
 - 2) Encode $B(x)$.
- } No bits shared!

Decompose:

① First decide $\delta(|B(x)|)$,

② Fetch $B(x) \Rightarrow$ get value of x in BAND.

$$S(x) = \underbrace{\delta(|B(x)|)}_{1 + 2 \cdot \log_2 L} \cdot B(x)$$

$$\text{Ex: } S(5) = \delta(|B(x)|) \cdot B(x)$$

$$= 011 \quad 101$$

DATA
AT COUNTER
 $\delta(|B(x)|)$

CARDINALITY
A.P. # Bits
used to represent

More space than Δ -encoding?

Let's actually check!

$$|S(x)| = \underbrace{L}_{B(x)} + \underbrace{2 \cdot \log_2 L}_{\delta(|B(x)|)} + 1 = |Enc(x)|$$

$$|S(x)| = \log_2 x + 2 \log_2 \log_2 x + 1$$

And we know that:

$$S(x) = \overbrace{2 \cdot \log_2 x}$$

$\Rightarrow S\text{-Encoder} \geq D\text{-Encoder}$,
 letter
 than
 when X no longer!

OPTIMAL DISSECTION FOR S -ENCODER.

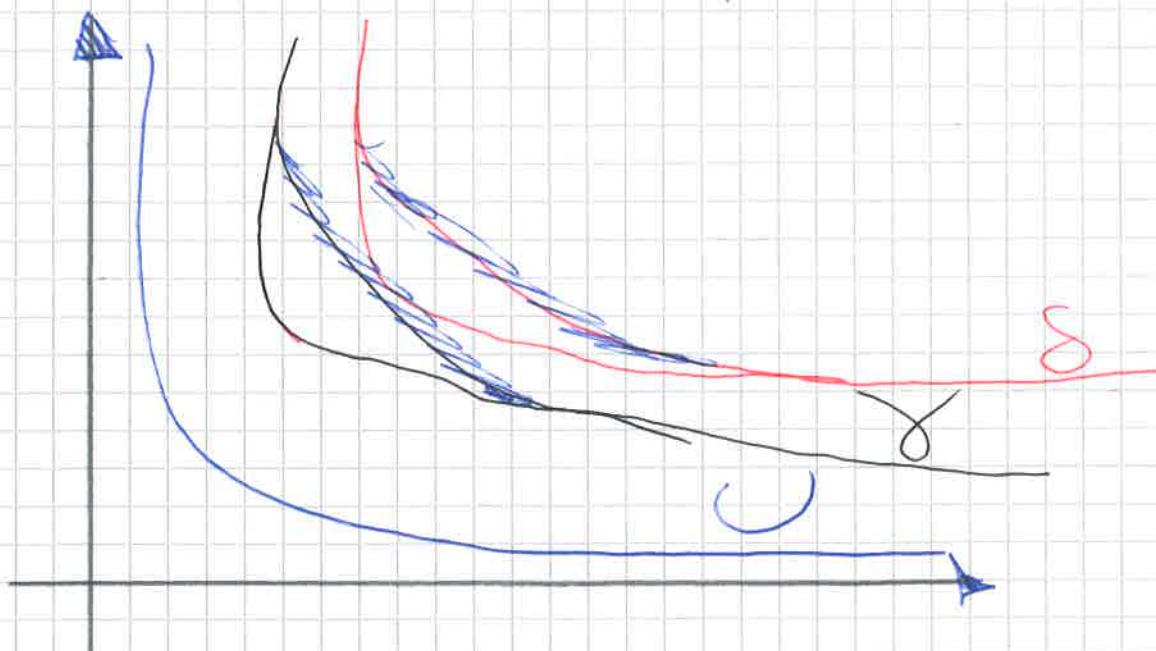
$$P_{\text{err}} | \text{Enc}(x) | = \log_2 \frac{1}{P\{x\}}$$

$\rightarrow \text{SOPSI}$

$$P\{x\} = \frac{1}{2 \cdot x / (\log x)^2}$$

$\Rightarrow 1 + O(N)$ from length of optimal binary code.

$\Rightarrow D$ and S -codes are UNIVERSAL and EQUIVARIANT when set S spread around 0.
 — slow DECODING (BUT SIMPLE)



OPTIMAL DISTRIBUTION FOR S-ENTROPY.

$$|E_{NC}(x)| = \log_2 \frac{1}{P\{x\}}$$

$$\log_2 x + 2 \log_2 \log_2 x + 1 = \log_2 \frac{1}{P\{x\}}$$

$$\log_2 x + \log_2 (\log_2 x)^2 + \log_2 2 = \log_2 \frac{1}{P\{x\}}$$

$$\cancel{\log_2(2x) + \log_2 (\log x)^2} = \log \frac{1}{P\{x\}}$$

~~$$\log (2x \cdot (\log x)^2) = \log \frac{1}{P\{x\}}$$~~

$$P\{x\} = \frac{1}{2x \cdot (\log x)^2}$$

RICE CODING:

Used when a set of integers are well-concentrated around a value $\neq 0$.

GOOD COMPRESSION RATIO & DECODING SPEEDS

PARAMETRIC ENCODING, based on parameter K

$$R_K(x) = \begin{cases} q = \left\lfloor \frac{x-1}{2^K} \right\rfloor & \text{QUOTIENT} \\ r = x - 2^K \cdot q - 1 & \begin{array}{l} \text{(represented in)} \\ \text{base } 2^K \end{array} \end{cases}$$

$$R_K(x) = \boxed{U(q)} \quad \boxed{r}$$

$q+1$ BITS K BITS

(q "0s" + 1 "1")

$q = \text{Quotient}$: Stored via $U(q)$ in a VARIABLE-LENGTH representation of $q+1$ BITS.

$$q < 2^K$$

$r = \text{REST}$: Stored via K BITS.

$$r \in [0, 2^K], \quad r < 2^K$$

PARAMETER K 's choice

Pick K in a way such that:

2^K is concentrated at the mean of the elements by S .

(The more concentrated are elements around 2^K , the better the encoding.)

Ex.: $R_4(83)$, where $K=4$.

$$R_4(83) = \left\{ \begin{array}{l} q = \frac{82}{2^4} = \frac{82}{16} = 5 \\ r = 83 - 2^4 \cdot 5 - 1 = 83 - 165 + 1 = 2 \end{array} \right.$$

$\underbrace{(1/6)}_{R_4(83)}$ $\underbrace{t \text{ in } k \text{ bits}}_{0010}$

$$|R_K(x)| = q + 1 + K \text{ BITS}$$

OPTIMAL DISTRIBUTION FOR RICE CODING

If values to be encoded follow a

GEOMETRIC DISTRIBUTION with parameter

(p) then we have:

$$P\{X\} = (1-p)^{x-1} \cdot p$$

Try to find the OPT. DISTRIBUTION for RICE CODING.

18 GEOMETRIC DISTRIBUTION specified:

$$P\{X\} = (1-p)^{x-1} \cdot p$$

$$\Rightarrow \lg 2^k \approx \frac{\ln 2}{p} \approx 0.69 \cdot \text{Avs/s} \text{ OPTIMAL}$$

→ This binning scheme can be used for encoding numbers ~~not~~ not only at the beginning of the distribution!

SPACE: $\left\lfloor \frac{x-1}{2^k} + 1 + k \right\rfloor$ Bits.

P-for DELTA ENCODING:

• EXTREMELY FAST DECOMPRESSION & small size in the compressed output, whenever:

REQUIREMENT:

S_n's values follow a GAUSSIAN DISTRIBUTION

PARAMETERS: (base, b)

~~base b = 3 13 39 37 33~~

IDEA: Force numbers repeating to ~~be encoded with few bits~~

(Need a coded sequence of integers to identify them!)

IDEA:

S =	BASE	12	15	16
BASE =	0	2	3	6
GAP-ENCODING	0	2	3	3

Pick a **Base** (usually the first number in the sequence, such that all numbers are non-negative).

$b = 2 = \# \text{BITS to encode } \rightarrow \text{the resulting number.}$

$n_b = \left\{ \begin{array}{l} n \text{ encoded with } b \text{ bits} \\ n \text{ expanded with } 2^b - 1 \text{ bits} \end{array} \right.$

$$\text{if } n \in \{0, 1, \dots, 2^b - 1\}$$

$$\text{if } n > 2^b - 1$$

\Rightarrow Put in a separate list to represent b bits.

$$0, 1, 1 + 2^{b-1}, \dots, 2^b - 1$$

EXAMPLE.

$S = 1 \quad 3 \quad 13 \quad 29 \quad 31 \quad 37$

S	BASE	0	2	12	28	30	36
GAP-ENCODING	0	2	1	16	2	6	
bin:	00	10	11	11	10	11	

(except)

\Rightarrow All values are encodable in fixed length b bits \Rightarrow

$b+6$ bits \Rightarrow DESCRIBED

N.B.: Escaped characters are encoded in a separate list using a fixed-size representation of $\binom{N}{b}$ bits (where $N = b$). They will be retrieved later on from a "special" data structure for escaped characters. How to pick b ? (greater b)

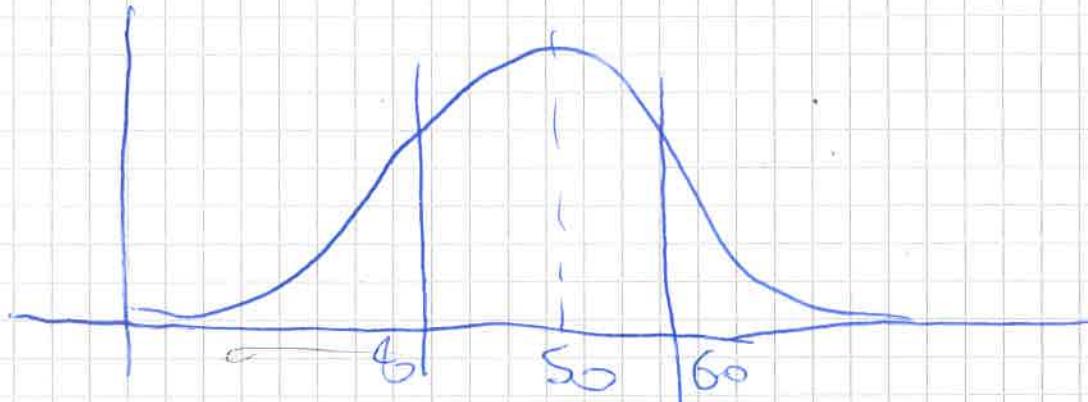
(1) METHOD: Trade-off between SPACE ALLOCATION and SPACE SAVING. Factor ENCODING to consider more repetitions smaller b .

(2) METHOD: Pick b in a way that 99% of values in S are $\leq 2^b$.

APPLICATIONS:

Virtually used in DATABASES!

Ex:



To keep an "array of signs":

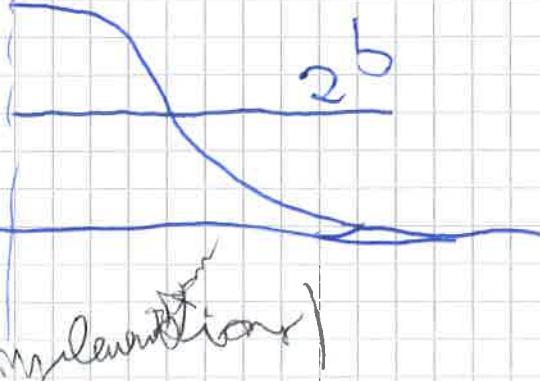
ADVANTAGE:

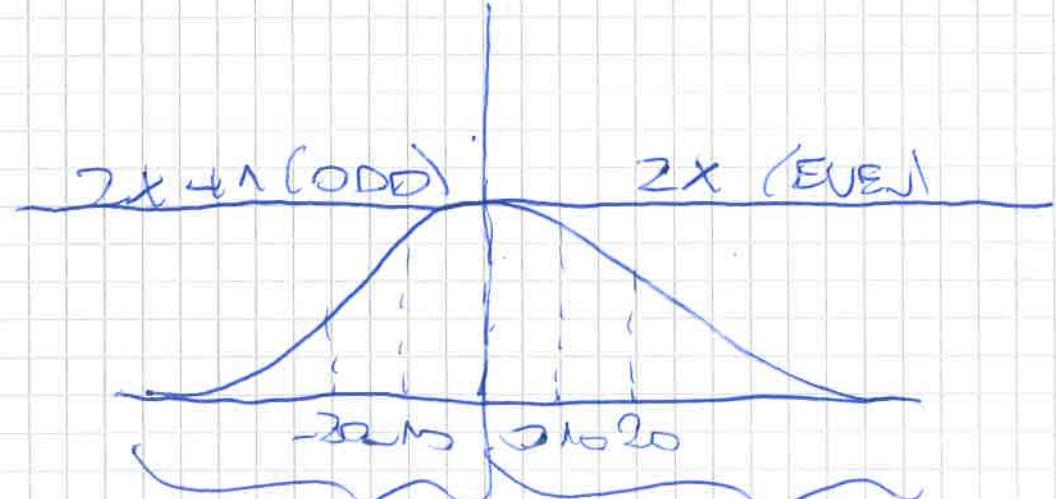
Extremely fast for decoding, as bits are word-aligned.

(Can also work

much faster.

IMPLEMENTATIONS

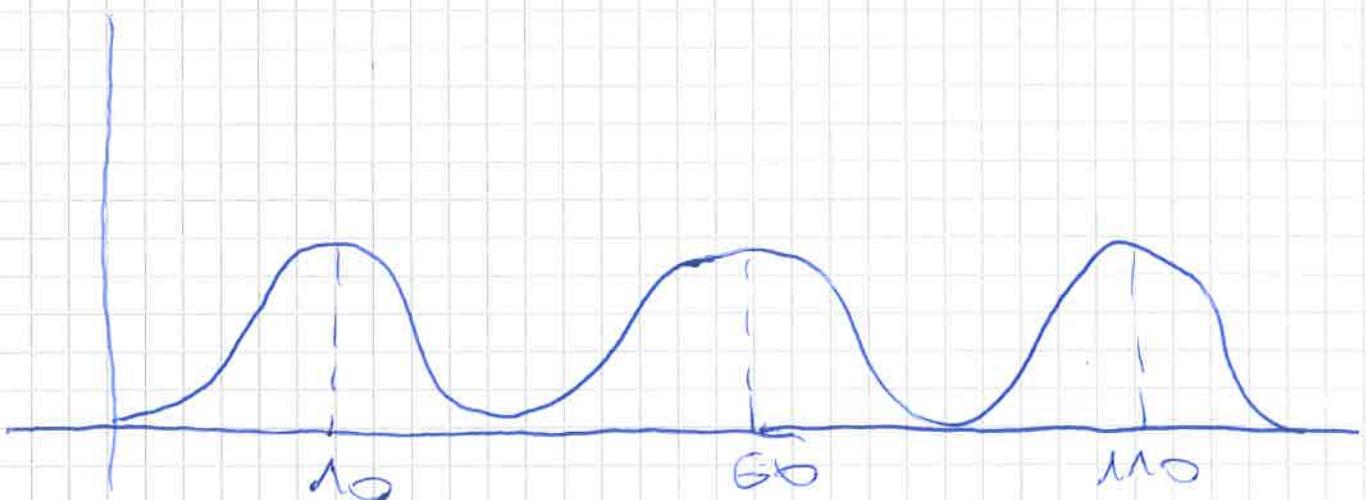




NEGATIVE NUMBERS
 by
 divided as ODD

POSITIVE NUMBERS
 by
 divided as EVEN

These don't work for a MULTI-MOD
 DISTRIBUTION.



ENCODING EXERCISES (in class)

D-CODING:

$$S = 1, 5, 10, 13, 15 \Delta S$$

$$\Delta S = 14 \quad 5 \quad 3 \quad 2 \quad S$$

2 GAP ~~100101~~

00101

Encoding

$$\delta(\Delta S) = 1 | 00100 | 00101, 011 | 010 | V$$

S-CODING: $\delta(\{S(x)\})$

$$S(4) = \underbrace{\delta(\beta)}_{100} = 011100$$

RKE-CODING:

$$R_3(\Delta) \xrightarrow{\text{divide}} R_3(\Delta S), k=3.$$

$$R_k(\Delta) = \left\lfloor \frac{\text{UNary}(\alpha)}{q+1 \text{ bits}} \right\rfloor \underbrace{k \text{ bits}}_{\text{K bits}}$$

q = QUOTIENT, r = REST in a UNARY STREAM.

r = REST, α = K BITS

$$\Delta S = 1 \quad 4 \quad 5 \quad 3 \quad 2 \quad S$$

$$\begin{aligned} & \begin{array}{c|c|c|c|c|c} 1 & 001 & 1 & 100 & 1 & 101 \\ \hline q+1 & \overbrace{1001}^{K \text{ bits}} & & & & \end{array} \\ & \alpha = \left\lfloor \frac{X=1}{2^k} \right\rfloor = \cancel{00100} = 0 \quad 1101 \end{aligned}$$

\Rightarrow For all encoding types \Rightarrow Say which distribution is optimal.

PFor Delta Coding.

$$S = \{1, 5, 10, 13, 15, 20\}$$

$$S_{\text{base}} = \{-9, -5, 0, 3, 5, 10\}$$

$$\text{Positive} = \{19, 11, 0, 6, 10, 20\}$$

$$\text{BIN}(x) = 111, 111, 000, 110, 111, 111$$

$$\text{ESCAPED SET} = \{19, 11, 10, 20\}$$

$$\text{base} = 10$$

$$b = 3$$

NB: Numbers are made positive according to

$$X = \begin{cases} -2X + 1 & \text{if } X < 0 \Rightarrow \text{ODD!} \\ 2X & \text{if } X > 0 \Rightarrow \text{EVEN!} \\ 0 & \text{if } X = 0 \end{cases}$$

In the $\text{BIN}(x)$ representation, we make use of 3 bits \Rightarrow Hence.

- We represent X up to $2^b - 1$ (z)

- We escape $X > 2^b - 1$ with:

$$\underbrace{111}_{5 \text{ bits}} = 2^b - 1 \text{ value for ESCAPE.}$$

\Rightarrow Escaped characters will be simply read from the list of escaped characters.

VARIABLE-BYTE CODE (S.C) - DECODE CODE.

Close & codes that trade space with
succinctness \Rightarrow One instantiation thereof in
the VARIABLE-BYTE CODE, which uses a sequence
of bytes to represent an integer X .
 \Rightarrow SIGNIFICANT DECODING SPEED.
REMARKABLE

IDEA: Take a number X and represent it in
BINARY $\Rightarrow B(X) \Rightarrow$ Partition $B(X)$ into groups
of 7 bits, padded at the beginning by
one BIT PADDING \rightarrow 0 : STOP (last group)
1 : CONTINUE (non-last group)

Decoder: Keep scanning the BYTE SEQUENCE
till we find a byte that is < 128 . (10000000)

EXAMPLE:

(e.g. Till we find
a STOPPED group)

$$2^7 = 128$$

INPUT: X

$B(X) \Rightarrow$ PARTITION into groups of 7 bits
(Add ~~padding~~ using 0's)
~~as padding is needed~~

EXAMPLE: Encode $X=2^{16}$ Group 0

$X=2^{16} \Rightarrow B(X) = 1000000000000000$

Group 2 Group 1

Group 2! Continue Bit \rightarrow stopping bit
00000100 | 00000000 00000000 = Encoding with
VARIABLE-BYTE CODE.

For any number x , we will have at least
 1 Group.

$$2^7 = 128$$

SPACE = At least 8 bits (7 for representation
 plus +1 for padding)

On average, + bits are wasted

Remember:

1 byte = 8 bits (given that x has 32 bits)

2 can represent $2^8 = 256$

SPACE in Bytes = $\frac{|\mathcal{B}(x)|}{7}$ (because one byte takes 8 bits
 & 1 bit is used for padding)

SPACE in Bits = $8 \cdot \frac{|\mathcal{B}(x)|}{7}$

N.B.: Drawback: When encoding very small digits, we waste lots of leading 0s for the representation!

Ex: 5
 $\Rightarrow 00000101$ (WASTED)

OPTIMAL DISTRIBUTION FOR VARIABLE BYTES CODE:

$$\underbrace{|\mathcal{E}_{\text{NC}}(x)|}_{\text{VARIABLE BYTES CODE}} = \log_2 \left(\frac{1}{P_{\{x\}}} \right)$$

$$8 \cdot \frac{|\mathcal{B}(x)|}{7} = \log_2 \left(\frac{1}{P_{\{x\}}} \right)$$

$$8 \cdot \log_2(x) = 7 \cdot \log_2 \left(\frac{1}{P_{\{x\}}} \right)$$

$$8 \cdot \log_2(x) = 7 \cdot \log_2 \left(\frac{1}{P_{\{x\}}} \right)$$

$$\log_{\text{O}_2}(x^8) = \log_{\text{O}_2} \left(\frac{1}{P\{X\}} \right)^7$$

$$x^8 = \left(\frac{1}{P\{X\}} \right)^7$$

$$x^8 = \frac{1^7}{(P\{X\})^7}$$

$$(P\{X\})^7 = \frac{1}{x^8}$$

$$P\{X\} = \sqrt[7]{\frac{1}{x^8}}$$

ADVANTAGE of VARIABLE-BYTE CODING:

GROUP-BASED decoding / Bit-by-Bit

IS very fast!

Based on this bit, performance
of the 2 actions?



- ① CONTINUE Decoding the next group

"CONTINUES"

CONFIGURATIONS

$$|C| = C = 2^7 = 128$$

- ② STOP Decoding the next group.

"STOPPS"

CONFIGURATIONS

$$|S| = S = 2^7 = 128$$

S-C-DENSE CODES:

$$\boxed{S+C=256} = \text{Overall \# possible configurations.}$$

$$= 2^8$$

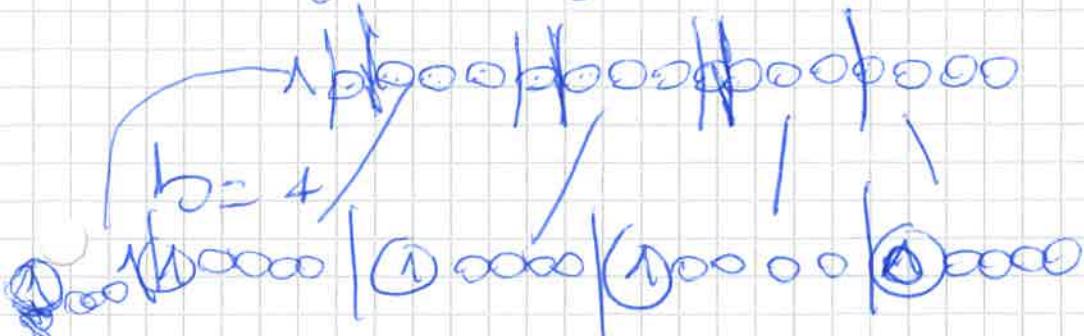
However, can we change the # bits used to divide the groups? (No longer int)

~~Ex:~~



Ex: Do not expand into groups of 7 bits, but

into groups of 8 bits.



\Rightarrow For a ~~GENERIC~~ GENERIC Group of size

2^b we will have

$$S + G = 2^b \Rightarrow \text{Based on the word size!}$$

#STOPPERS
CONFIGURATIONS
in b bits.

#CONTINUED CONFIGURATIONS
in b bits

\Rightarrow The value to pick for S and G depends on the DISTRIBUTION of integers to be encoded.

How do we handle configurations?



CONFIGURATIONS' COMBINATIONS:

(i.e., what are valid combinations of configurations?)

#BUTTED

#GAPS

CONFIGURATIONS

#CONFIGURATIONS (among bits)

1



S

2



$C \cdot S$

3



$C^2 \cdot S$

For n groups/bytes:

$\Rightarrow n$ bytes can encode integers.

$$\# \text{CONFIGURATIONS} = C^{n-1} \cdot S / C^{n-1} \cdot S$$

~~Note: n Groups~~

Ex: $b=3 \Rightarrow \# \text{CONFIGURATIONS} = 2^b = S + C$

$2^3 = S + C \Rightarrow$ We can have encode values

from $[0, 15]$

$3+1$ $5+1$ $7+1$

min min min

$S=C \Rightarrow$ Encode values up to S .
write it.

VALUES: $S=C=4$ $S=6, C=2$

		STUFFED GROUPS	
0	4	0 0 0 0	0 0 0
1	VALUE	0 0 1	0 0 1
2	ENCODE	0 1 0	0 1 0
3		1 1 1	0 1 1
4		1 0 0 0 0 0	1 0 0
5		1 0 0 0 0 1	1 0 1
6		1 0 0 0 1 0	1 1 0 0 0 0
7		1 0 0 0 1 1	1 1 0 0 0 1
8		1 0 0 1 0 0	1 1 0 0 1 0
9		1 0 0 1 0 1	1 1 0 0 1 1
10		1 0 0 1 1 0	1 1 0 1 0 0
11		1 0 0 1 1 1	1 1 0 1 0 1
12		1 1 0 0 0 0	1 1 1 0 0 0
13		1 1 0 0 0 1	1 1 1 0 0 1
14		1 1 0 0 1 0	1 1 1 0 1 0
15		1 1 0 0 1 1	1 1 1 0 1 1

100

Difference between ~~$S=5, C=4$~~ and

$$\textcircled{2} S=6, C=2 \quad b = 2^6 = 64$$

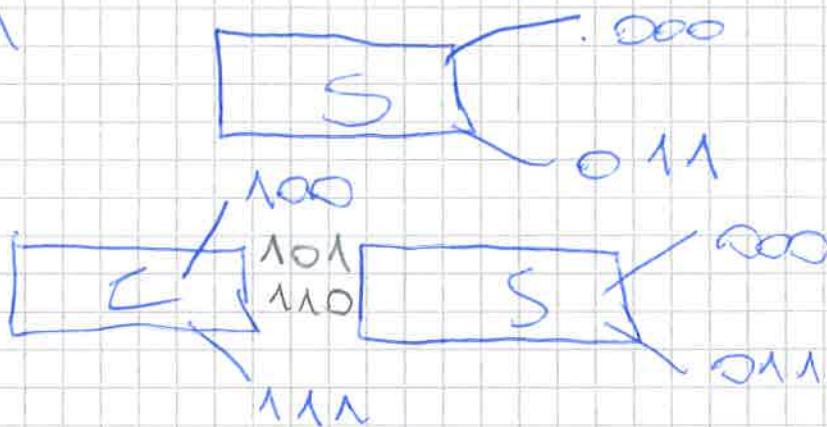
$$b=2 \Rightarrow b-1=2^6-1=63$$

In later cases,

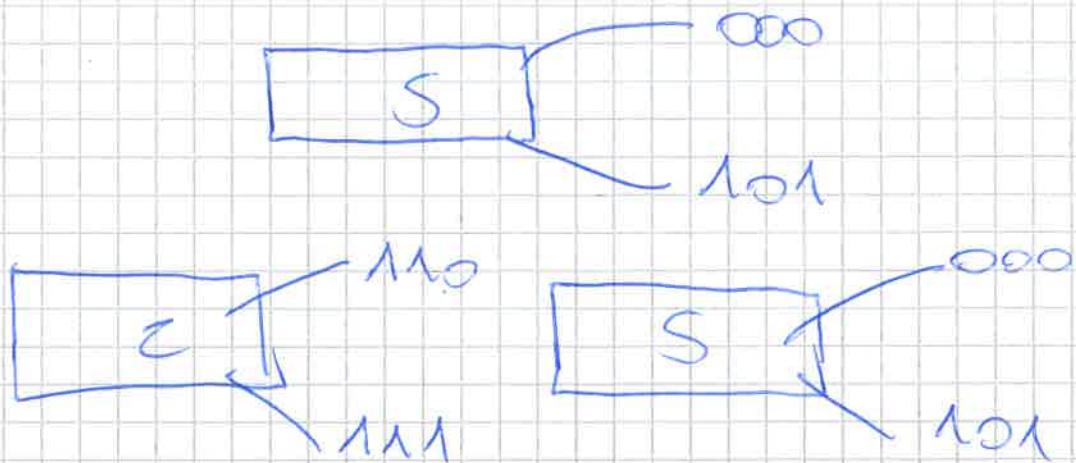
we use 2 words of 3 bits or 6 bits to encode
all the ~~small~~ integers.

In case (1): $S=4, C=4$, we encode the first
4 values with DNF word.

(1)



In case (2): $S=6, C=2$, we decode the first
6 values with ONE word, and the following ~~any~~
~~with~~ 2 values, with



This can be made advantageous, based on
the distribution of $\{1, 15\}$.
 S is small (smaller in size).

~~(S, c)~~ DENSE CODE - CLASS

Exercises

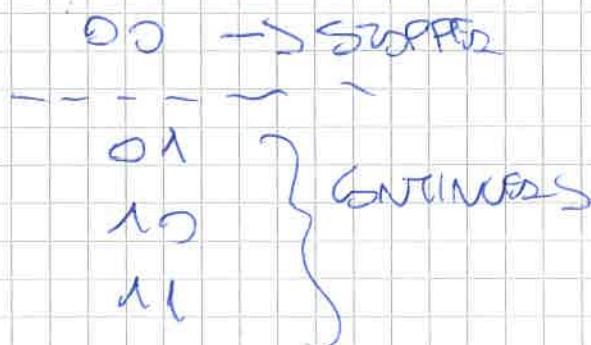
$$S=1, C=3$$

(1, 3) - Code

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

[Must be a power of 2]
2 Bits!
per group.

Configurations



• CONSTRUCT the FIRST 15 codewords

Blocks	USES	RANGE & COMBINATIONS	INDEX	ENCODING
$S=1$	\boxed{S}	1 [0]	①	00
C	\boxed{S}	$S=1$ $C=3$ $*3[1-3]$	1 2 3	01 00 10 00 11 00
C	C	S	4-12	01 01 00 01 10 00 00 11 00 10 01 00 10 10 00 00 11 00 10 11 00 11 01 00 11 10 00 11 11 00
C	C	C	13-15	00 01 01 00 00 01 10 00 00 10 10 00 00 11 11 00
C	C	C	16-17	00 01 01 10 00 01 11 00
C	C	C	18-19	00 10 01 00 00 10 11 00
C	C	C	20-21	00 11 01 00 00 11 11 00

~~We keep going down as long as:~~

~~MIN(n) > i, where i is
the partition position of a .~~

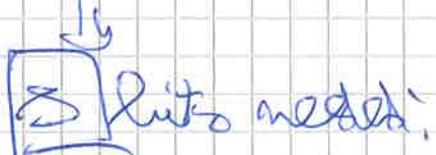
~~→ At the end we return~~

$$i - \text{MIN}(n), \text{ where } n \text{ is the last index we visited.}$$

~~→ MIN(n), where n is the last index we visited.~~

EX: How many bits do I need to represent the number 15?

→ Need 4 containers & 1 topper



CONFIGS:

$$S=C=4 \quad C=2; S=6$$

~~WORST~~

S=C=4	C=2; S=6	WORST
4	6	<input type="checkbox"/> S
16	12	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> S
64	24	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> S
		6 111 101

\Rightarrow S and C can be chosen by minimizing CONNECT FUNCTION.

APPROACH TO ENCODE a number:

$$C=S=4; b=3$$

$$X=8$$

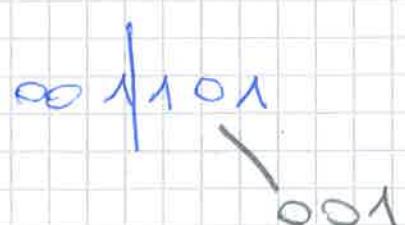
$$B(x) = \cancel{000} \ 001 \mid 000$$



~~$B(x) = 00001000 \quad C=2; S=6$~~

$$X=13$$

$$B(x) =$$



INTERPOLATIVE CODE.

WHAT?

Integer-encoding technique ideally used when there is a CLUSTERED occurrence of integers [in small ranges] (EX: Store Posting lists \rightarrow ZIP codes in search engines)

INPUT: Increasing sequence S (sorted)

$$S_{i-1} < S_i < S_{i+1}$$

Ex: # PREVIOUS ELEMENTS # FOLLOWING ELEMENTS to the MAX.

$$S = \textcircled{1} 2 3 5 7 \textcircled{11} 15 18 19 20 \textcircled{21}$$

OFFSET = 1 2 3 + 5 6 7 8 9 10 11 12

We repeat the fact that (21) is the max. number and that the numbers are increasing & DISTINCT.

IDEA: We want to encode (9) w.r.t. its lower bound & its upper bound.

$$\boxed{\text{Upper Bound}} = \text{MAX} - \# \text{ FOLLOWING ELEMENTS to the MAX.}$$

Assuming S is increasing by 1 at every step till the end

$$\begin{aligned} &= 21 - 6 \\ &= 15 \end{aligned}$$

$$\boxed{\text{Lower Bound}} \# \text{ ELEMENTS PRECEDING to the MIN}$$

$$= 6.$$

$$\# \text{ CONFIGURATIONS} = 9 - 6 = 3$$

$$\# \text{BITS} = \log_2 (15 - 6 + 1) = \log_2 10 = 4$$

All Possibilities of Offset
Upper Bound / Lower Bound

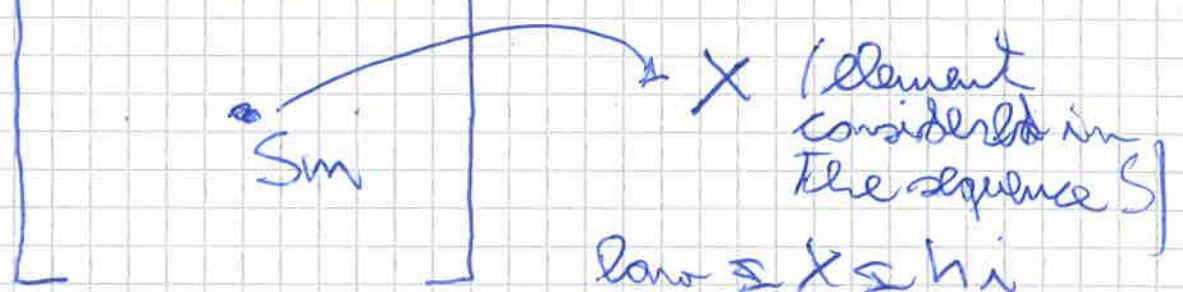
$$\# \text{BITS} = \lceil \log_2 (\text{upper} - \text{lower} + 1) \rceil$$

Very efficient encoding w.r.t. space.

Idea: Each number \rightarrow in S is encoded with as few bits as possible.

We here obtain an IMPACT Encoding when considering the sequence of increasing numbers.

low = Lower Bound hi = Upper Bound.



BINARY CODE: $\langle X, \text{low}, \text{hi} \rangle$

We consider the ~~Binary Encoding of~~ $X - \text{low}$ in $\lceil \log_2 (\text{hi} - \text{low} + 1) \rceil$ bits with

WHAT: $\boxed{X - \text{low}}$ as encoded in **BINARY**.

SPACE: $\lceil \log_2 (\text{hi} - \text{low} + 1) \rceil = n$
(real number)

~~$S_{\text{seq}} = S_1, S_2, \dots, S_m, \dots, S_{n-1}, S_n$~~

RECURSIVE ENTHALV PROCEDURE:

Let's consider our INCREASING & DISTINCT sequence S:

$$S = 1 \ 2 \ 3 \ 5 \ 7 \ 9 \ 11 \ 15 \ 18 \ 19 \ 20 \ 21$$

$$Q = 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11 \ 12$$

At each RECURSIVE STEP, we know:

- l = left - index of $S_{l,r}$; $low = \min(S_{l,r})$
- r = right - index of $S_{l,r}$; $high = \max(S_{l,r})$
- w = length of $S_{l,r}$

Where $S_{l,r}$ is the sequence within a RECURSIVE CALL.

INITIALIZATION:

$$n = |S|; l = 1; low = S_1; hi = S_n; r = n$$

Stops in the compressed file for ~~decoding~~
allowing the reader to decide them.

Maintenance:

In a recursive call, we operate similarly to
Binary Search. Namely:

$$m = \lfloor \frac{l+r}{2} \rfloor$$

& Allows us to find S_m

INVARIANT: In each Recursive call, we have:

- $\forall s_i \in S_{l,r}: \boxed{\begin{array}{l} s_i > low \\ s_i < hi \end{array}}$

Recursion's Idea:

We must encode the whole sequence (S) via a binary-matching procedure. Now,

Interpolative Coding ($S[1, r], \text{low}, \text{hi}$):

if $r < 1$: // we have reached the end of a string
return \emptyset

end if

if $l == r$, then:

return Binary Code ($S[l]$, low, hi)

end if

compute $m = \lfloor \frac{l+r}{2} \rfloor$

$S[l] - \text{low}$

$\overset{l}{\underset{m}{\text{---}}}$

~~Compute $A_1 = \text{BinaryCode}(S[m], \text{low}+m-1, \text{hi}-l+m)$~~
~~if encode $S[m]$,~~

$A_1 = \text{BinaryCode}(S[m], \text{low}+m-1, \text{hi}-l+m)$

$A_2 = \text{InterpolativeCoding}(S[1, m-1], \text{low}, S[m]-1)$
|| go left of $S[m]$

$A_3 = \text{InterpolativeCoding}(S[m+1, r], S[m])$
|| go right of $S[m]$ (hi)

return concatenation of $A_1 \cdot A_2 \cdot A_3$;

VISUALIZED: LEFT,



InterpolativeCoding($S[1, m-1], \text{low}, S[m-1]$)
 $S[1, r]$ low hi

VISUALIZED: RIGHT



Interpolative Coding ($S[m+1], S[m+i], hi$)
 $S[1, r]$ low m m

NB: low and high are VALUES!

and r and l are POSITIONS!

The encoding of an integer S_i is not FIXED, but it depends on the distribution of other integers as well (Variable Encoding & ADAPTIVE).

NOT prefix-free.

⇒ EXPLOIT DISTINCTNESS & INDEPENDENCE of the elements in the sequence.

INTERPOLATIVE CODING VS SHANNON'S THEOREM

In certain cases, we are even able to do better than Shannon's what Shannon's theorem dictates.

⇒ No Bits needed for coding numbers!

Ex: 1, 2, ..., 8.

After putting the first tuple, moving limits.

EXAMPLE: Interpolative Coding ($S'[1, 12], 1, 21$) from left

$S = \underline{1} \underline{2} \underline{3} \underline{5} \underline{7} \underline{9} \underline{11} \underline{15} \underline{18} \underline{19} \underline{20} \underline{21}$

~~Root node~~ $m = \lfloor \frac{1+12}{2} \rfloor = 6$

$A_1 = \text{BinaryCode}(9, 1 + 6 - 1, 21 - 12 + 6)$

$\text{Binary Code } (9, 6, 15) \Rightarrow 9 - 6 = \boxed{3} \text{ FMT}$

LEFT:

$\text{IntCoding}(S'[1, 5], 1, 7)$ $\text{IntCoding}(S'[7, 12], 11, 21)$

$1 \ 2 \ 3 \ 5 \ 7$

$11 \ 15 \ 18 \ 19 \ 20 \ 21$

$m = \lfloor \frac{1+5}{2} \rfloor = 3$

$\text{BinaryCode}(3, 1 + 3 - 1, \dots)$

$\begin{matrix} \\ 0 \\ \diagup \\ \end{matrix}$

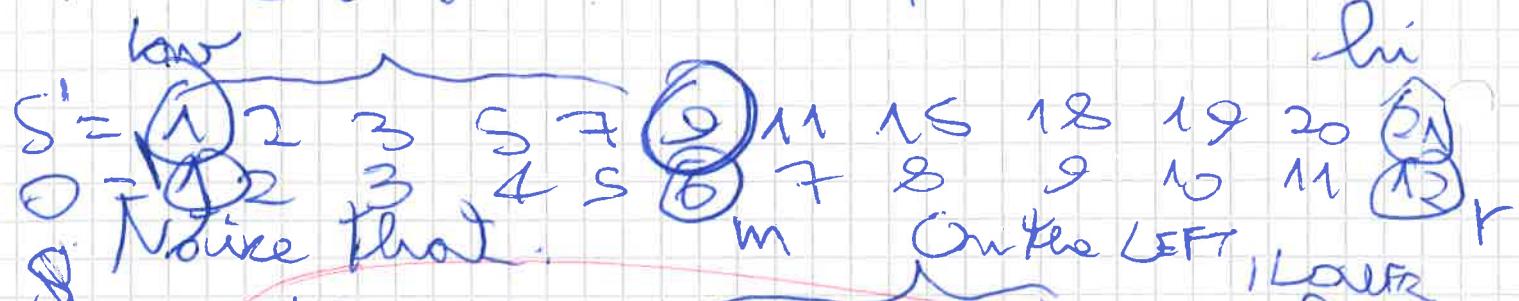
$\text{IntCoding}(S'[4, 5], 5, 7)$

$\text{IntCoding}(S'[1, 2], 1, 2)$

$m = \lfloor \frac{1+2}{2} \rfloor = 1$

BinaryCode

How to encode? \Rightarrow Binary Encoding



LEFT: $S'[m] \geq \text{low} + m - 1$
 $\geq 1 + 6 - 1$

We have in fact in S', t :

$m - 1$ distinct values in the left-hand side
and smallest one is $\geq \text{low}$.

RIGHT: $S[m] \leq \text{hi} - (r - m)$ (UPPER BOUND)
 $S[m] \geq \text{hi} - r + m$

\Rightarrow We hence ~~do not~~ Do NOT encode $S[m]$,

but $S[m] - (\text{low} + m - 1)$ Known Lower Bound.

Encoder with: $\log_2(S[m] - (\text{low} + m - 1))$ Bits

We hence know that:

$$\text{high} + m - 1 \leq S[m] \leq \text{low} + m - 1$$

$$\text{low} + m - 1 \leq S[m] \leq \text{high} + m - 1$$

Drawback: To decode one element, needs to know its lowest element.

INTERPOLATIVE Coding - CLASS EXERCISE

$$S = [2, 5, 8, 10, \underbrace{11, 12, 13}_{\text{1 2 3 4 5 6 7}}]$$

$$l = 1; r = 7; (n=7); \text{low} = 2; \text{hi} = 13$$

$$m = \lfloor \frac{l+r}{2} \rfloor = 4 \Rightarrow S[m] = 10$$

~~Binary Encode~~ \rightarrow Binary Encode $(\text{low} + m - 1, \text{hi} - r + m)$

$$\text{Binary Encode } (5, 1) \Rightarrow S[m] = 5 = (5)$$

~~How~~ in the #BITS needed to represent the range!

$\Rightarrow S$ is represented in 3 bits.

Result \rightarrow ~~EMIT~~: 101

~~FFFF~~ $\rightarrow S[m]$

$$l = 5; r = 7; \text{low} = 11; \text{hi} = 13; n = 3$$

~~Binary Encode~~

$$m = \lfloor \frac{l+r}{2} \rfloor = 6$$

Binary Encode $(12, 12)$

$$\Rightarrow \text{Encode } \overset{S[m]}{\cancel{12}} - 12 = 0 \Rightarrow \text{No Bits}$$

~~LEFT~~ \rightarrow ~~S[m]~~ \rightarrow ~~emits~~

~~Binary Encode~~

$$l = 1; r = 3; \text{low} = 2; \text{hi} = 8; n = 3$$

$$m = \lfloor \frac{l+r}{2} \rfloor = 2 \Rightarrow S[m] = 5$$

$$\text{Binary Encode } (3, 7) \Rightarrow \text{Encode } S[m] - 3 \Rightarrow 5 - 3 = 2$$

Ende (2) mit 3 Lits

EMV \approx 10

ELIAS-FANO ENCODING:

INTERPOLATIVE Coding!

WHAT?

Code that does NOT depend on the distribution of integers to be encoded, allowing for INDEXING through compressed data structure.

WHEN?
RANDOM

By Efficiently ACCESS the ~~original~~ encoded integers in a compressed way

(Ex: Used: Store inverted lists ~~in~~ search engines & adjacency lists of large graphs).

BAD: Space is a CONCERN ~~for~~ or integers have clusters

REQUIREMENTS:

(S) is a monotonically increasing SEQUENCE
(like for INTERPOLATIVE Coding) & all DISTINCT numbers.

$$S_{i-1} < S_i < S_{i+1}$$

Ex: 1 3 5 7 12 15 18 20

Each INTEGER is represented in BINARY:

~~b = Flag 2 0 = # bits~~

~~intvec() bin() & Mat. integers~~

~~1
3
5
7
12~~

~~FUNCTIONAL~~ EXPLAINED.

Let's take $S = 1 \ 3 \ 9 \ 12$

Each integer is represented in BINARY as:

#BITS per INTEGER = $b = \lceil \log_2 (U) \rceil$, where
 $(U = \max S_i)$ (i.e. max. integer in S)

S_i	BN/S_i	$b = \lceil \log_2 12 \rceil = 4$
1	0001	$b = \lceil \log_2 U \rceil$
3	0011	$l + w = b$
9	1001	$w = b - l$ or
12	1100	$l = b - w$

Now let's consider a positive integer l and a positive integer U to partition the BINARY REPRESENTATION into 2 BLOCKS.

$$l = \lceil \log_2 \frac{U}{n} \rceil = 2 = \# \text{BITS less significant}$$

$n = 15$

$$w = \lceil \log_2 (U+n) - l \rceil = 4 - 2 = 2 = \# \text{BITS more significant}$$

We call the two representations respectively:

- $L(S_i)$: Block of the l least significant bits (lower ones)
- $H(S_i)$: Block of the $b-l$ most significant bits.

And we will surely have that.

$$|H(S_i)| + |L(S_i)| = b \text{ bits}$$

$$|L(S_i)| = l \cdot n \text{ bits} = \log_2(S_i)$$

$$|H(S_i)| = n + 2^w \text{ bits} \Rightarrow \begin{matrix} \text{SPACE FOR} \\ \text{MOST} \\ \text{SIGNIFICANT} \\ \text{BITS} \end{matrix}$$

OVERALL # ONES OVERALL # ZEROES

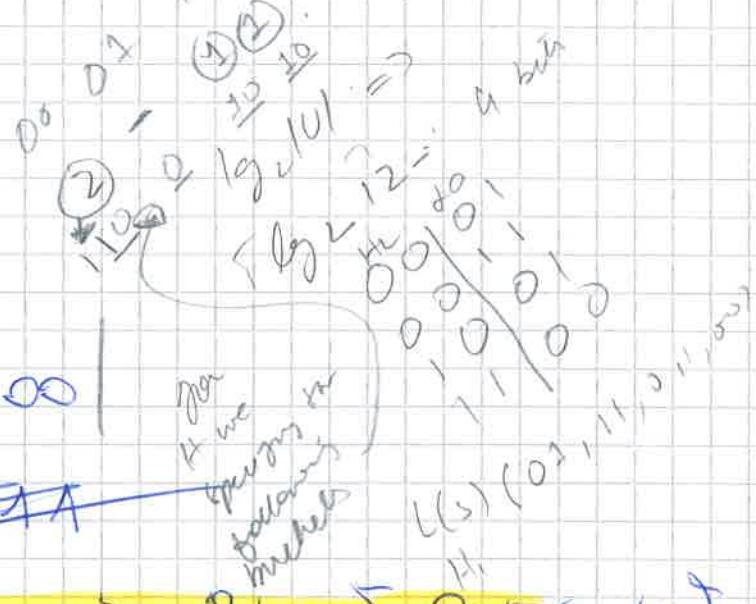
EXAMPLE:

S_i	$B(S_i)$
00	0001
01	0011
10	1001
11	1100

$$\boxed{L(S_i)} = 01\ 11\ 01\ 00$$

$$\cancel{H(S_i) = 00\ 00\ 10\ 11}$$

$$H(S_i) = 10$$



We group the most significant bits into ~~BUCKETS~~ ~~BUCKETS~~ ~~BUCKETS~~ ~~BUCKETS~~ ~~BUCKETS~~ ~~BUCKETS~~ ~~BUCKETS~~ ~~BUCKETS~~ BUCKETS, then ~~Count~~ ~~Count~~ ~~Count~~ ~~Count~~ ~~Count~~ ~~Count~~ ~~Count~~ represent each bucket by the Count of groups in that BUCKET in ~~BINARIES~~ ~~UNARY~~.

For H , we identify the following BUCKETS

BUCKETS	00	01	10	11
Groups	2	0	1	1
Groups' Count (UNARY)	110	0	10	10

$$\boxed{= H(S_i)}$$

~~$H(S_{it})$ has 2^w possible configurations.~~

$$\text{SPACE} = |H(S_{it})| + |H(S_{it})|$$

$$= \cancel{n \cdot 2^w}$$

$$n \cdot |L(S_{it})| = n \cdot (\log_2 \frac{U}{n})$$

The FINAL ENCODING is given by:

$$\boxed{H(S_{it}) \oplus L(S_{it})}$$

CONCATENATION

$$= \cancel{n \cdot 2^w} \quad \cancel{\log_2 \frac{U}{n}} \quad \cancel{\log_2 \frac{U}{n}}$$

$$|L(S_{it})| = n \cdot (\log_2 \frac{U}{n}) = \begin{matrix} \text{Space of each} \\ \text{integer} \end{matrix} \quad \begin{matrix} \text{bits} \\ \text{2 least significant} \end{matrix}$$

$$|H(S_{it})| = \cancel{n} + \cancel{2^w} = \begin{matrix} \text{All possible configurations} \\ \text{(Bucketed) with } w \text{ bits.} \end{matrix}$$

$$\text{SPACE} \approx |L(S_{it})| + |H(S_{it})|$$

$$= \cancel{n \cdot (\log_2 \frac{U}{n})} + \cancel{n + 2^w}$$

$$\text{Where } w = \log_2(U+1) - \log_2 \frac{U}{n}$$

$$\approx \log_2 U - \log_2 \frac{U}{n}$$

$$\approx \log_2 \frac{U \cdot n}{U}$$

$$\Rightarrow w = \log_2 n$$

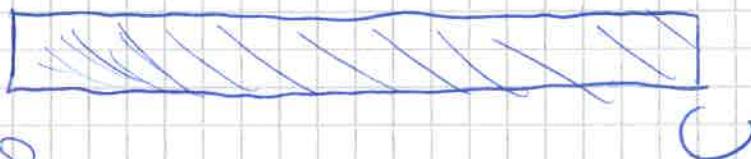
↪ Let's plug in this value of λ :

$$\begin{aligned}\text{SPACE}_{\text{IMBAL}} &\approx n \cdot \left(\log_2 \frac{\lambda}{n} \right) + n + 2 \log_2 n \\ &\approx n \cdot \left(\log_2 \frac{\lambda}{n} \right) + n + n \\ &= \boxed{n \cdot \left(\log_2 \frac{\lambda}{n} \right) + 2n} \\ &= n \cdot \left[\log_2 \frac{\lambda}{n} + 2 \right]\end{aligned}$$

Regardless of the distribution! 1/2 bits per integer

Very close to the optimal space!

LOWER BOUND for SPACE - OPTIMAL:



Let's consider the possible combinations of n elements into a space of size λ .

combinations =

$$\binom{\lambda}{n}$$

$$\text{SPACE for all combinations} = \log_2 \binom{\lambda}{n} \approx n \cdot \log_2 \frac{\lambda}{n}$$

$$\Rightarrow \boxed{n \cdot \log_2 \frac{\lambda}{n}} \leq n \cdot \left(\log_2 \frac{\lambda}{n} \right) + 2n$$

Lower Bound! \Rightarrow ELIAS-FANO

STIRLING'S FORMULA

IS Elias-Gono takes 2 bits more per integer (n) than the OPTIMAL LOWER BOUND (if integers are uniformly distributed.)

ELIAS-FAND

DECOMPRESSION: Easy & random access allowed; just by counting the # ZEROS.

ELIAS-FANO CODING (LESS EXERCISE)

$$S = \{N_1, S_1, N_2, N_3, N_5, D_2\}$$

<u>Si</u>	<u>B(Si)</u>
000 - [1]	00001
001 - [5]	00101
010 - [10]	01010
011 - [13]	01101
100 - [15]	01011
101 - [20]	10101
110 - [25] = $\log_2(10)$	10111
111 - [30] = $\log_2(11)$	10101

$$n = |S| = 6$$

3 4 5
8 16 32

$$l = \lceil \log_2 \frac{17}{n} \rceil = \lceil \log_2 \frac{20}{6} \rceil = \lceil \log_2 3.2 \rceil = 2$$

$$x = \log_2(0+n) - l = 5 - 2 = 3$$

$$e + v = b \quad \checkmark$$

$$L(Si) = \sigma_1 / (\sigma_1 + \sigma_2 + \sigma_3 + \sigma_4)$$

BUCKETS: ONE FOR EACH SERVICE THAT USES THE API

Graves' 1 1 1 2 1 1 0 0

Count:

FINAL ENCODING $W(S)OL(S)$

• **Voluptuous**

(a) 2

HUFFMAN CODING.

Part of the larger family of STATISTICAL CODING! \rightarrow It involves 2 phases:

- **MODELLING PHASE:**

Compute statistical properties of input & build a model.

- **CODING PHASE:**

Actually compresses the input sequence.

STATISTICAL CODING can be applied to:

- Sequences of ^{replicates} symbols ($TEXT$) drawn from an alphabet Σ / of finite size

\hookrightarrow **TEXT ENCODING**.

- Sequences of **INTEGERS** (numbers), drawn from an alphabet Σ / of integers

\hookrightarrow **INTEGER CODING**

- **GENOMES CODING**

\Rightarrow **OPTIMAL ENCODING** can be solved for any distribution via HUFFMAN CODING / at the cost of storing the PREAMBLE (n.o: some overhead).

If an OPTIMAL distribution by the data to encode is possible PREFIX-FREE, then we ~~can't~~ go for it \rightarrow to reduce the overhead and avoid complicated decoding.

Have to HUFFMAN-ENCODED i.e. function

Recall the definition of entropy H :

$$H = \sum_{\theta} P\{\theta\} \cdot \log_2 \frac{1}{P\{\theta\}},$$

where θ is a symbol in Σ alphabet.

Code's length: (AVERAGE CODEWORD'S LENGTH)

$$L_C = \sum_{\theta} P\{\theta\} \cdot \text{len_code}(\theta)$$

SHANNON'S THEOREM:

$$\sum_{\theta \in \Sigma} P\{\theta\} = 1$$

If for every symbol we have a code, then a LOWER BOUND holds!

$L_C \geq H$ holds, where:

$$0 \leq H \leq \log_2 |\Sigma|$$

HUFFMAN FUNKTION:
ENCODING

Huffman Coding is based on a GREEDY LOCAL CHOICE, that eventually leads to an OPTIMAL situation \Rightarrow BUILD HUFFMAN TREE!

INPUT: INITIALLY: We have a set of symbols θ taken from Σ , each one with $P\{\theta\}$ of occurring \Rightarrow These symbols constitute the LEAVES of our tree.

SYMBOLS & A CORRESPONDING PROB. DISTRIBUTION.

a b c d e f
0.05 0.1 0.15 0.3 0.25 0.15

CANDIDATE NODES

LEAVES

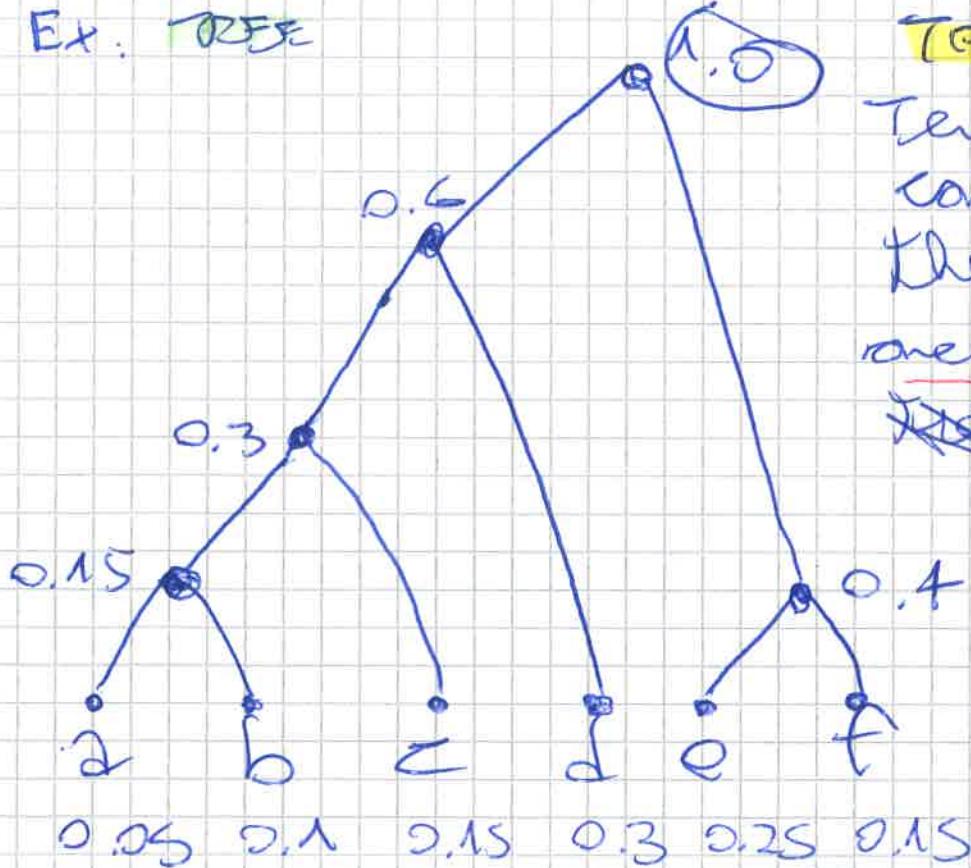
from CANDIDATE NODES

- * REPEAT at every step :
- ① Pick 2 leaves / nodes with the lowest probability & merge them into a new node \Rightarrow add to candidate set
 - ② Sum their probability P_{TOT} } into the newly created node & remove merged nodes from the candidate set

[NB: In case of a TIE among the sum of nodes' probability, just pick one pair \Rightarrow The tree won't change, yet OPTIMALITY is preserved]

HUFFMAN

Ex.: TREE



TERMINATION:

The algorithm continues until there is only one node ~~left~~ (the root) with $P_{\text{TOT}}=1$.

⇒ How to get the HUFFMAN CODE based on the HUFFMAN TREE?

(2) Get the HUFFMAN CODE based on the HUFFMAN TREE, how? By LABELED HUFFMAN TREE all we assign BINARY LABELS (0 or 1) to the TREE EDGES of the HUFFMAN TREE.

↳ Generally:

- ↳ LEFT EDGE of an internal node ⇒ 0
- ↳ RIGHT EDGE of an internal node ⇒ 1

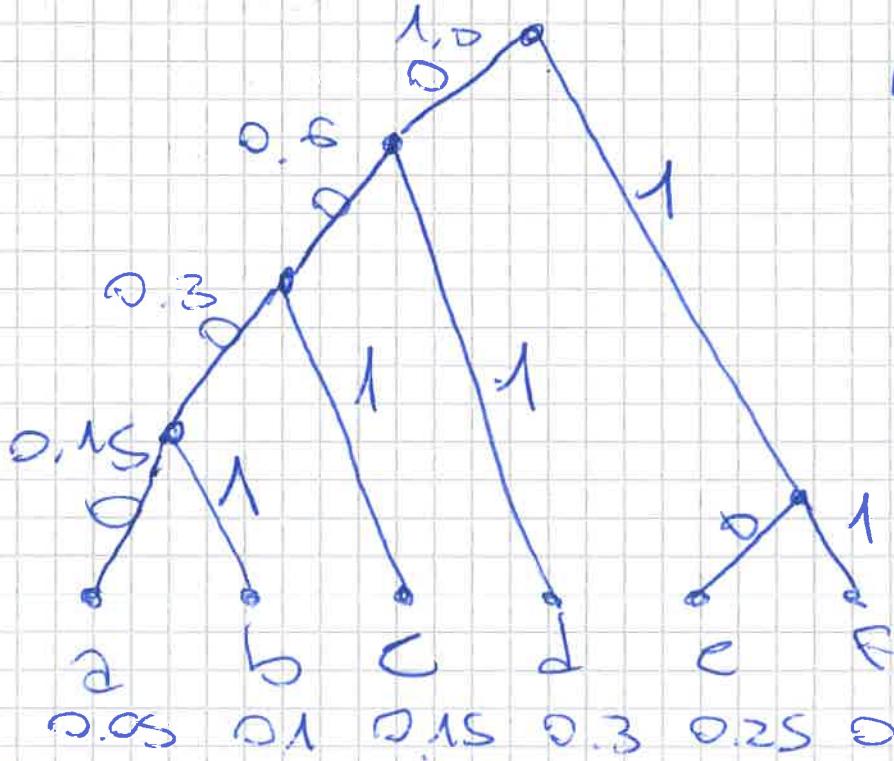
$$\# \text{ INTERNAL NODES in a HUFFMAN TREE} = n - 1$$

$$\# \text{ POSSIBLE HUFFMAN CODES in a HUFFMAN TREE} = 2^{\sum_{i=1}^{n-1} 1} = 2^{n-1}$$

[2 CHOICES: 0 or 1, in each edge, for each internal node]

↳ Based on a HUFFMAN TREE, we can then obtain the ~~HUFFMAN CODE WORD~~ LABELED HUFFMAN CODE WORD for a symbol Ø by:
- Taking the BINARY LABELS encountered on the downward path connecting the root to the leaf Ø.

Ex. Build LABBED HUFFMAN TREE & create a HUFFMAN CODEWORD.



LEFT EDGE = 0
RIGHT EDGE = 1

ORDER

RULE OF

~~greedy~~

statistical data coding:

The rarer a symbol, the more bits we use to encode it.

~~greedy~~

To encode it.

The more frequent a symbol, the fewer bits we use to encode it.

CODEWORD of $\theta = 2 = 0000$
 $b = 0001$
 $c = 001$
 $d = 01$
 $e = 10$
 $f = 11$

N.B. We notice that making a different choice in picking symbols with probability has an impact on the MAX. codeword length, which would impact on the size of the compression / decompression buffer.

1 IDEA to minimize the **MAXIMUM LENGTH**; we consider the OLDEST nodes

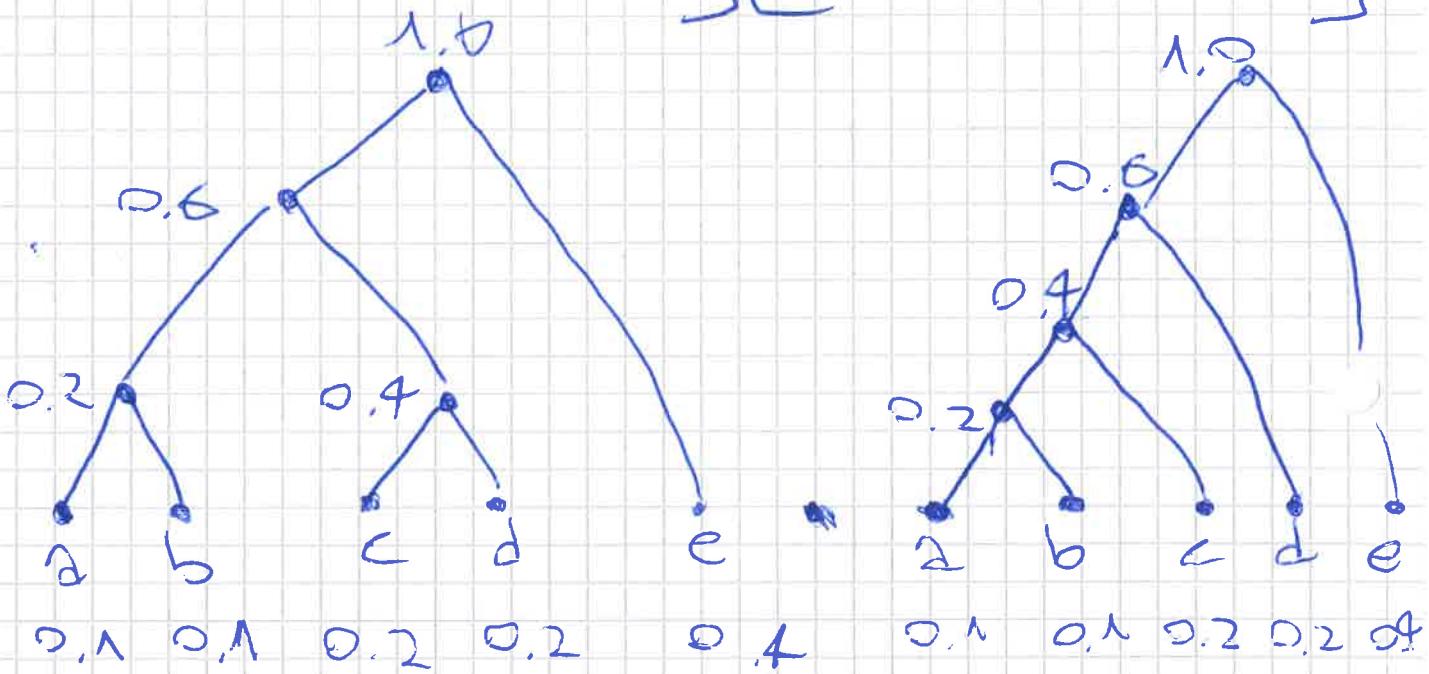
\Rightarrow OLDEST Nodes \rightarrow Leaves or internal nodes that have been merged farther in the past than the other nodes in the candidate set.

\Rightarrow The 2 queues:

- ↳ (1) QUEUE contains symbols sorted by increasing probability
- ↳ (2) QUEUE contains internal node in the order they are created by the Huffman algorithm.

TIE:

[Closest to QUEUE 1; then in QUEUE 2] [APPROACH - Queue - best] [APPROACH - Random]



Huffman is NOT the only optimal encoding algorithm \Rightarrow You can have more Huffman encoding algorithms that are optimal as well.

\Rightarrow Since Huffman is optimal, then it must be optimal as well.

COMPRESSED FILE, following HUFFMAN'S ALGORITHM

HUFFMAN APPLICATION:



0/1SII
(HUFFMAN TREE)
Code words &
symbols 0 & 1

Avg. HUFFMAN CODEWORD LENGTH:

L_H = Average Codeword length produced by Huffman

$$H \leq L_H \leq H+1 \quad = \text{OPTIMAL}$$

\Rightarrow The Avg. codeword length is hence bounded by the entropy $H = \sum_{i \in S} p_i \cdot \log_2 \frac{1}{p_i}$

Knowing the entropy H , we can hence estimate codeword's length

THEOREM - OPTIMALITY OF HUFFMAN CODE:

If C is a HUFFMAN CODE, then L_C is the shortest possible Avg. length among all the PREFIX-FREE codes C' , where,

$$L_C^H \leq L_{C'} \quad \begin{array}{l} \text{codewords of other} \\ \text{prefix-free code } C' \end{array}$$

$$L_C^H = \sum_{i \in S} p_i \cdot l(i)$$

PROOF: We notice that a prefix-free code can be seen as a BINARY TREE.

\Rightarrow We show the minimality of Avg. length path to prove HUFFMAN CODE'S OPTIMALITY.

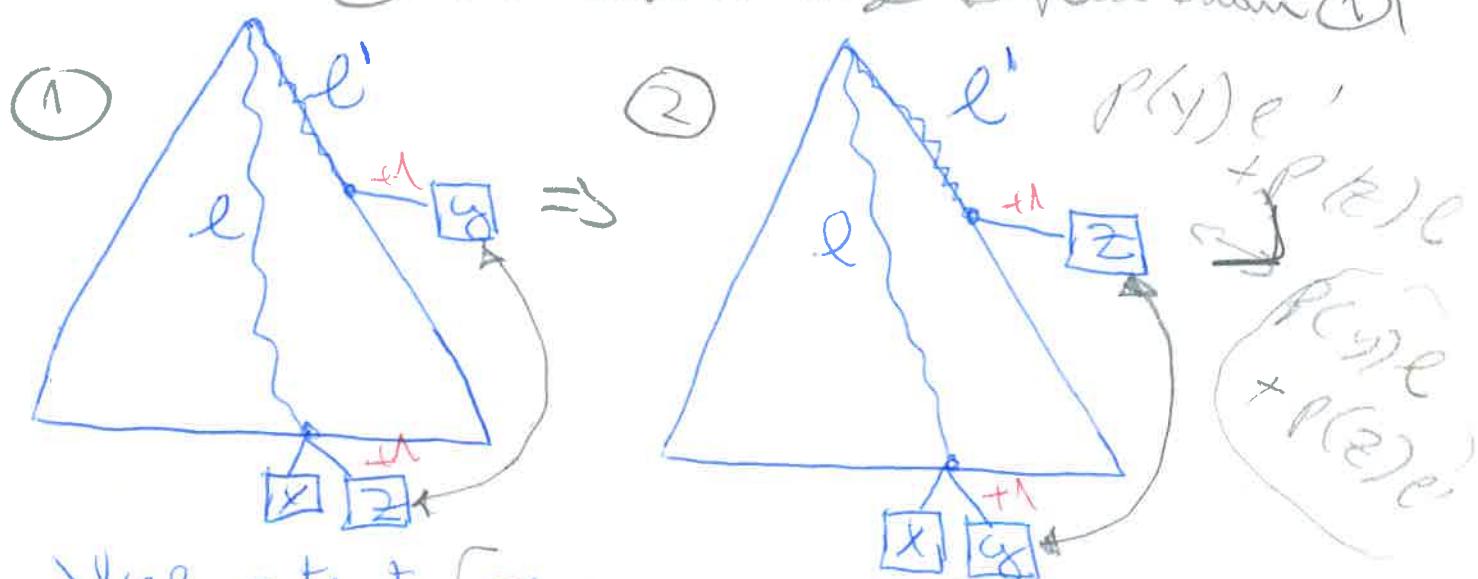
\Rightarrow how to prove:

If the theorem does not hold, then:

- A NOT minimum probability leaf (say ②) occurs at the deepest level of the binary tree. ① configuration
 - ~~②~~ ② can hence ~~not~~ be swapped with a MIN. probability leaf (say ①) to reduce the Ave. depth of the resulting binary tree.
- ② configuration

That \Rightarrow If configuration ② is less advantageous than configuration ①.

(i.e.: ② has flat average depth than ①)



We have that: $[P\{z\} \geq P\{y\}, P\{y\} \geq P\{x\}]$

- ① ~~$D_l + 1 \cdot P\{x\}$~~ + $(l + 1 \cdot P\{y\}) + (l' + 1 \cdot P\{z\}) + \dots$
- ② ~~$D_l + 1 \cdot P\{x\}$~~ + $(l + 1 \cdot P\{y\}) + (l' + 1 \cdot P\{z\}) + \dots$
- ③ ~~$D_l \cdot P\{z\} + P\{z\}$~~ + ~~$P\{y\}$~~ + $D_{l'} \cdot P\{y\}$ + $P\{y\}$
- ④ ~~$D_l \cdot P\{y\} + P\{y\}$~~ + $D_{l'} \cdot P\{z\}$ + ~~$P\{z\}$~~

~~PROOF BY INDUCTION~~ following

⇒ We now show that indeed:

Configuration (1) → Configuration (2)

$$l \cdot P\{z\} + l' \cdot P\{y\} > l \cdot P\{yz\} + l' \cdot P\{z\}$$

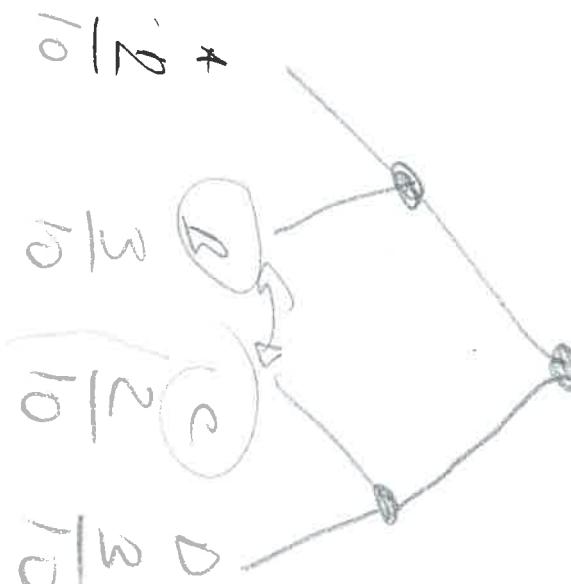
$$P\{yz\}/(l-l') > P\{yz\}/(l-l')$$

LEMMA \downarrow $\Rightarrow P\{yz\} > P\{yz\}$ Q.E.D!

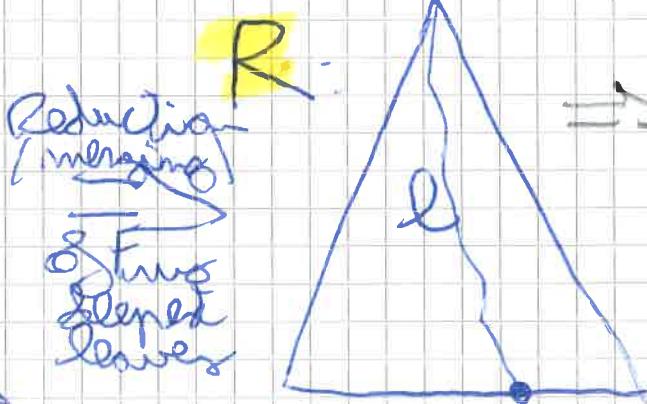
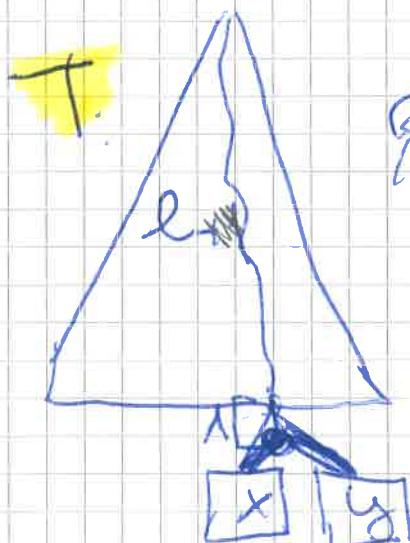
Let T be a BINARY TREE whose avg. depth is minimum among the binary trees with $|S|$ leaves.
 \Rightarrow Then, the two leaves with MIN. probability will be at the greatest depth of T and will be children of the same parent node.

PROOF: Consider an alphabet Σ with n symbols, and (X) and (Y) have the smallest Prob. (i.e.: they are Candidates with minimum P .)

$$\begin{aligned} P(X)(e) \\ P(Y)(e) \end{aligned}$$



The situation is hence the following.



\Rightarrow REDUCED TREE

$$P\{Z\} = P\{x_3\} + P\{y_3\}$$

L-T:

$$(l+n)(P\{x_3\} + P\{y_3\}) + \Delta \quad | \quad l(P\{x_3\} + P\{y_3\}) + \Delta$$

$$l \cdot (P\{x_3\} + P\{y_3\}) + (P\{x_3\} - P\{y_3\}) + \Delta$$

L-T

$$l \cdot (P\{x_3\} + P\{y_3\}) + \Delta$$

L-R

$$\Rightarrow L-T = L-R \text{ When? } \boxed{\text{LEMMA 2:}}$$

$$(L-T) = L-R + P\{x_3\} - P\{y_3\}$$

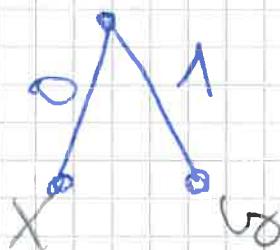
We now want to prove the OPTIMALITY of HUFFMAN CODE, by drawing that among other code, say Σ , would be longer than HUFFMAN CODING.

OPTIMALITY OF HUFFMAN CODING.

$$\forall \Sigma, L\Sigma \geq LH$$

PROOF By **INDUCTION** on n , The #symbols in Σ .

- **BASE CASE:** Take $n=2$ symbols.



\Rightarrow Obviously OPTIMAL, as it assigns one ~~the~~ single bit 0 to one symbol and ~~one~~ single bit 1 to the other symbol.

- **INDUCTIVE STEP:** Let's consider $n \geq 2$ symbols in Σ , and assume:

Assume:
HYPOTHESIS: Huffman code is OPTIMAL for $n-1$ symbols.

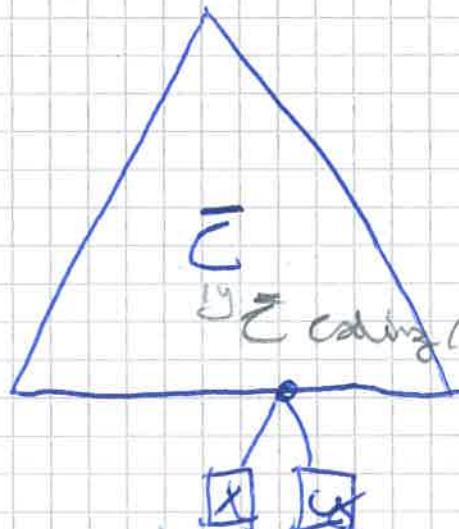
\rightarrow Take $|\Sigma|=n$ and let Σ be an optimal code for Σ and its underlying distribution.

We will want to show that:

$$L\Sigma = LH$$

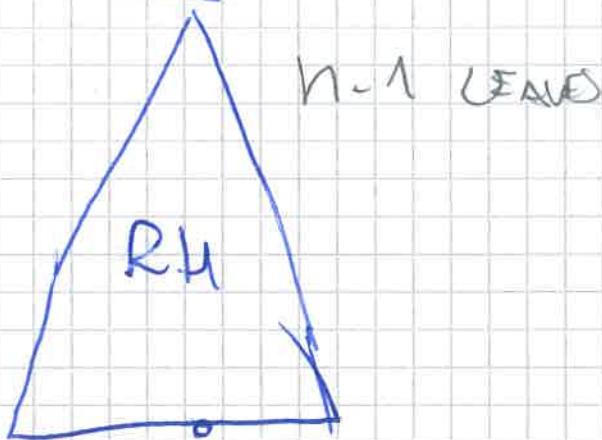
(i.e. Huffman is also OPTIMAL)
and this will have to
hold for (n) symbols.

Let's now consider two BINARY TREES with n leaves.



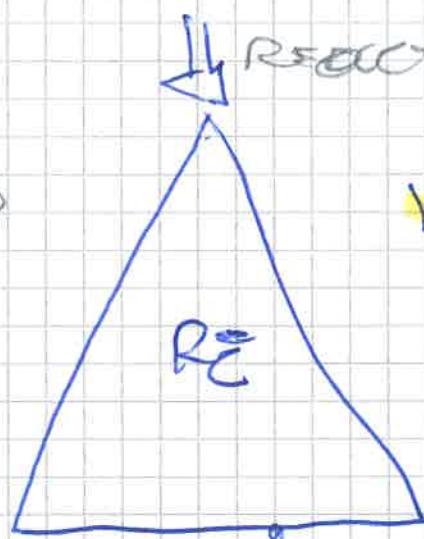
n leaves

↓ Reduction



$P_X + P_Y$

↓ Reduction



$n-1$ leaves

$P_X + P_Y$

Following a REDUCTION, R_H and R_E have respectively $n-1$ leaves.

⇒ Let's now apply our inductive hypothesis for R_H and R_E , which have $n-1$ leaves.

- By INDUCTIVE HYPOTHESIS, R_H is OPTIMAL for the reduced symbols (where X and Y have been merged). Hence:

$$AVG. \text{ depth} \text{ of an HUFFMAN CODE} \leq R_H \leq L_{R_E} \leq AVG. \text{ depth of an OPTIMAL CODE after the reduction}$$

We can now apply LEMMA 2 over

$$L_{RH} \leq \cancel{L_{R\bar{z}}}$$

Where we know that by LEMMA 2.

$$L_T = L_Z + (P\{x\} + P\{y\})$$

n lower n-1 lower

Hence, in the present situation:

$$L_H = L_{R\bar{z}} + (P\{x\} + P\{y\})$$

And by LEMMA 2 again:

$$L_{\bar{z}} = L_{R\bar{z}} + (P\{x\} + P\{y\})$$

$$\left\{ \begin{array}{l} L_{RH} = L_H - (P\{x\} + P\{y\}) \\ L_{R\bar{z}} = L_{\bar{z}} - (P\{x\} + P\{y\}) \end{array} \right.$$

SUBSTITUTING above we get:

$$\Rightarrow L_H - (P\{x\} + P\{y\}) \leq L_{\bar{z}} - (P\{x\} + P\{y\})$$

$$\Rightarrow L_H \leq L_{\bar{z}}$$

Q.E.D.

We have shown the optimality of HUFFMAN CODE.

LOSSLESS ENCODING WITH HUFFMAN.

We always pay some overhead to store the preamble ~~header~~, which contains our HUFFMAN TREE ~~header~~.

However, if our HUFFMAN TREE is small, the preamble is negligible.

$$H = -H + H + 1 \Rightarrow H \text{ overall} = 1 \text{ extra bit per symbol}$$

If we have a very large H , 1 bit per symbol is negligible.

Otherwise, if we don't have a very large H , 1 bit can have an impact.
Unknown word ~~a symbol with~~ ≤ 1 bit

\Rightarrow HUFFMAN & SHANNON'S INTUITION,

If you want to encode a very large text, say of infinite size, you can chunk it into K chunks of FINITE SIZE(K).

$$T = \frac{1}{K} + \frac{1}{K} + \frac{1}{K} + \frac{1}{K} + \dots$$

$\Rightarrow \sum^1 = \{ \text{blocks of } k \text{ symbols} \}$

+ extended alphabet of size \sum^k , whose symbols are substrings of k -symbols.

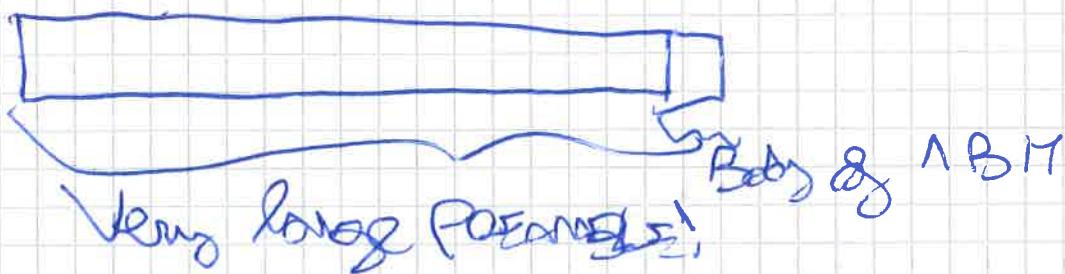
⇒ If we use Huffman on the ~~blocks~~ blocks, the extra bit is lost!

⇒ The larger is the block, the smaller is the penalty I would pay (asymptotically optimal)

$$H < L_H < H + \frac{1}{K}$$

However, this would increase the size of the Huffman tree, which contains the PREAMBLE!

RESORTING SITUATION, as $T \rightarrow D$.



My in practice, impractical!

N.B.: K is fixed, even for an infinite tail!

⇒ Still, most compressors make use of Huffman in some way (MPEG1/SACD).
(Trade off between LEAVES / symbols and BLOCKS).

$$H = \frac{1}{\log_2 |\Sigma|}$$

though,

CANONICAL HUFFMAN CODE

Huffman code is strongly limited by 2 sets:

- It has to store the structure of the tree in the preamble (only if Alphabet Σ is large!)
- Decoding is slow, as we need to traverse the tree from root to leaf. \rightarrow At every edge of the path we may have a CACHE MISS.

NEW IDEA: ~~Implement~~ implement a variant of the Huffman Code (Canonical Huffman) that alleviates such problems \rightarrow VERY FAST DECODE

No NEED to navigate ~~in~~ SMALL MEMORY
The HUFFMAN TREE all the time! CACHE

\rightarrow We need to make one ~~extra data structures~~
& operate in the following way:

① Compute codeword length $L(\theta)$ for each symbol $\theta \in \Sigma$ according to classical Huffman's Algorithm.

② Construct the array ~~num~~^{NCW}, which stores in array position $num[l]$ the #SYMBOLS having Huffman codeword of l -bits.

③ Construct the array ~~Symbol~~^{MATRIX}, which stores in the entry $symbol[l]$ the first l -symbols having Huffman codeword of l -bits.

(4) Construct the array ~~symbol~~ fcv, which stores in the entry $fcv[l]$ the first codeword of ~~symbol~~ all symbols included with ℓ bits according to measure (4)

(5) Assign consecutive codewords to the symbols in symbol[l], starting from the codeword $fcv[l]$.

• STORAGE:

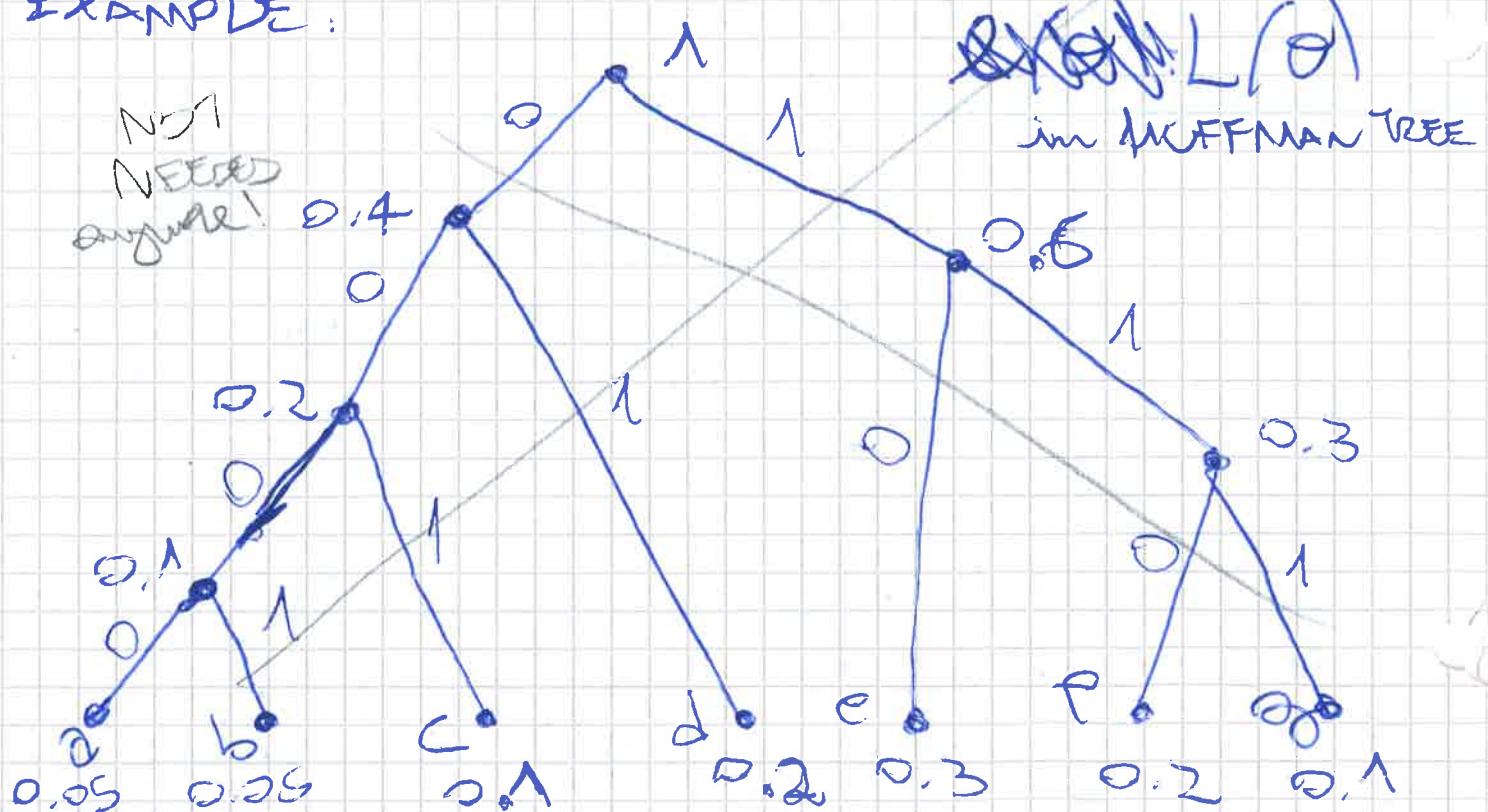
ADVANTAGE: We no longer need to store the tree-structure via pointers, but just store ~~fc~~ and ~~symbol~~ arrays. \rightarrow Other

• DECODE:

Data structures are discarded

We no longer need to traverse the whole Huffman tree, but we only operate on two available arrays (at most 2 code mixes per symbol instead of all labels along the path)

EXAMPLE:



construct array num / how

1	2	3	4
0	2	3	2

→ \forall level

$\in [0, \text{maxLevel}]$

construct array symbol.

array
levels

1	
2	d e
3	c f g
4	a b

We can fit millions
of symbols into cache!

construct binary F_{level}

according to
procedure (4)

1	2	3	4
2	3	1	0

$F_{\text{level}}[l]$ is the value
of the code word consisting
of l bits

$$F_{\text{level}}[i] = F_{\text{level}}[i-1] + \frac{\text{num}[i]}{2}$$

? i
to odd

DECOMPRESS

⇒ We can now convert ~~decompress~~ via:

$V = \text{next_bit}(); l = 1$

while ($V < \delta_{\text{level}}[l]$):

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

$l++;$

return symbol[l, $V - \delta_{\text{level}}[l]$]

IDEA: When FZW reaches the right of a code word to decide, then we actually decide.
[Reaching the right level!]

NB: No code misses occur ~~at the~~ within the while, but only in the return.

⇒ improves ENCODING & DECOMPRESSION SPEED largely.

Is $P_1 \cdot P_2 = \frac{1}{2+2}$, is it still worth using running counts? No, just go for a D-CODE, as it is much faster to apply.

To build ~~F~~ FZW:

$$fzw[\max] = 0$$

$$\text{for } l = \max - 1; l \geq 1; l--$$

$$fzw[l] = \frac{fzw[l+1] + new[l+1]}{2}$$

SYMBOL TABLE:

Level \ 0 1 2 3
 \ 0 1 2 3

1 D C A

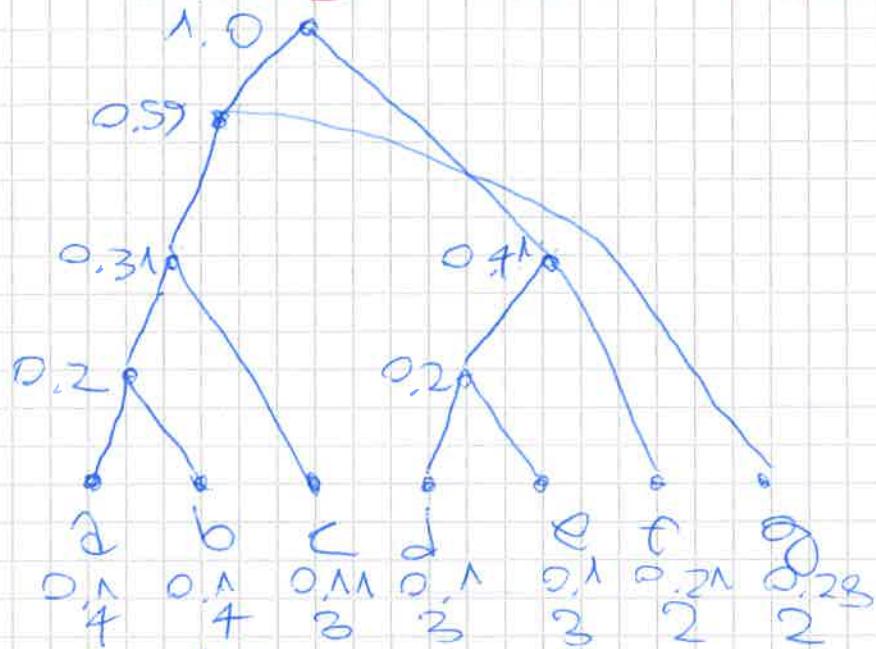
2 E

3 F

+

L

CANONICAL



HUFFMAN CODE - COMPRESSION

Build the tree:

1	2	3	4
0	2	3	2

SYMBOL TABLE

level	0	1	2	3
1				
2	f	g		
3	c	d	e	
4	a	b		

$$CH(\theta) = FCW[\ell(\theta)] + \text{offset}(\theta) \quad \text{initial BMS}$$

$$CH(g) = 2 + 1 = 3 \text{ in } 2 \text{ bits} = (11)_2$$

$$CH(f) = 2 + 0 = 2 \text{ in } 2 \text{ bits} = (10)_2 \Rightarrow FCW[2]_2$$

$$CH(e) = 1 + 2 = 3 \text{ in } 3 \text{ bits} = (011)_2$$

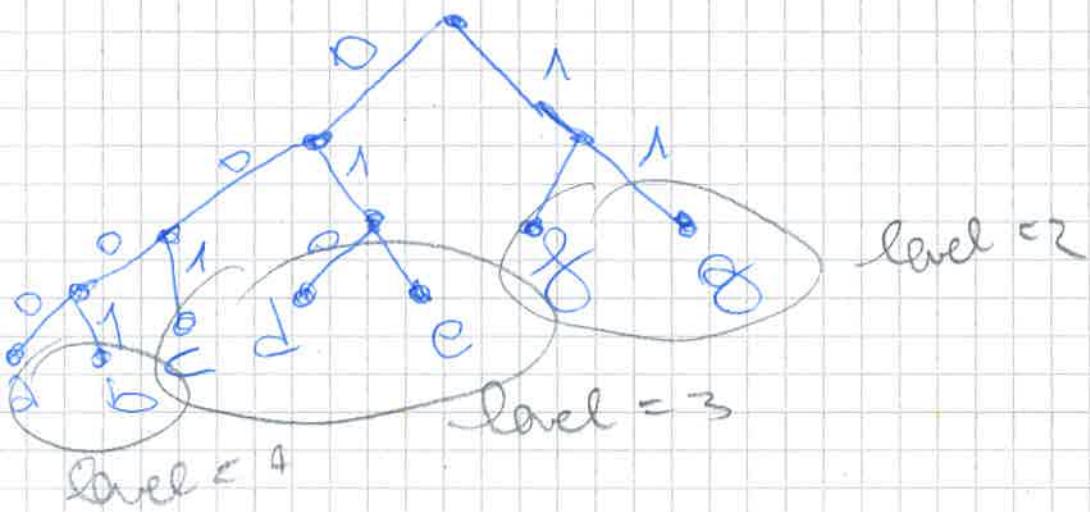
$$CH(d) = 1 + 1 = 2 \text{ in } 3 \text{ bits} = (010)_2$$

$$CH(c) = 0 + 1 = 1 \text{ in } 3 \text{ bits} = (001)_2 = FCW[2]_2$$

$$CH(b) = 0 + 1 = 1 \text{ in } 4 \text{ bits} = (0001)_2$$

$$CH(a) = 0 + 0 = 0 \text{ in } 4 \text{ bits} = (0000)_2 = FCW[4]_2$$

→ Alternatively, we could build a
CANONICAL HUFFMAN TREE based on
their codeword.



~~Bloom FILTER - CLASS EXERCISE~~

$$k_{opt} = \frac{m}{n} \ln 2 = \lceil 16 \ln 2 \rceil$$

$$E_{opt} = (0.6185) \frac{m}{n}$$

Error

~~Compute size of Bloom FILTER (m) to guarantee error is 2^{-20} bits~~

Canonical HUFFMAN - DECODING

$$v = next_bit(i)$$

$$l = 1;$$

while ($v < FOL[l]$)

$$v = 2v + next_bit();$$

$i++;$ last

return symbol[$N - 8L[l]$]

~~QUESTION WE HAVE : S = D1~~

$$FOL = \begin{array}{c} ① \\ \hline 2 | 3 | 4 \\ \hline 2 | 1 | 1 | 0 \end{array}$$

plan

SYMBOL TABLE		1	2	3	4
offset	symbol	D	C	A	B
1	E				
2					
3	F				

1) $l = 1$

$S = 01$

$V = 0$

$V \leftarrow \text{Fst}[1]$

$0 < 2 \Rightarrow$ Go into the while loop

$V = \cancel{0} 0 + 1$

$I = 2$

2) $l = 2$

$V = 1$

$V \leftarrow \text{Fcv}[2]$

$1 < 1$ ↗

→ Return symbol^e $\cancel{S}, V - \delta e[l]$

symbol[2; 1 - $\delta e[2]$]

symbol[2; 0] = D (first symbol)
 $\delta e[2]$

~~Idea: Keep examining bits until we find~~

~~$\delta e[l]$ - $\cancel{S}[l]$~~

~~Then, we return \cancel{S}~~

~~Symbol. $\text{SYMB}[0, \cancel{S}[i] - \delta e[l]]$~~

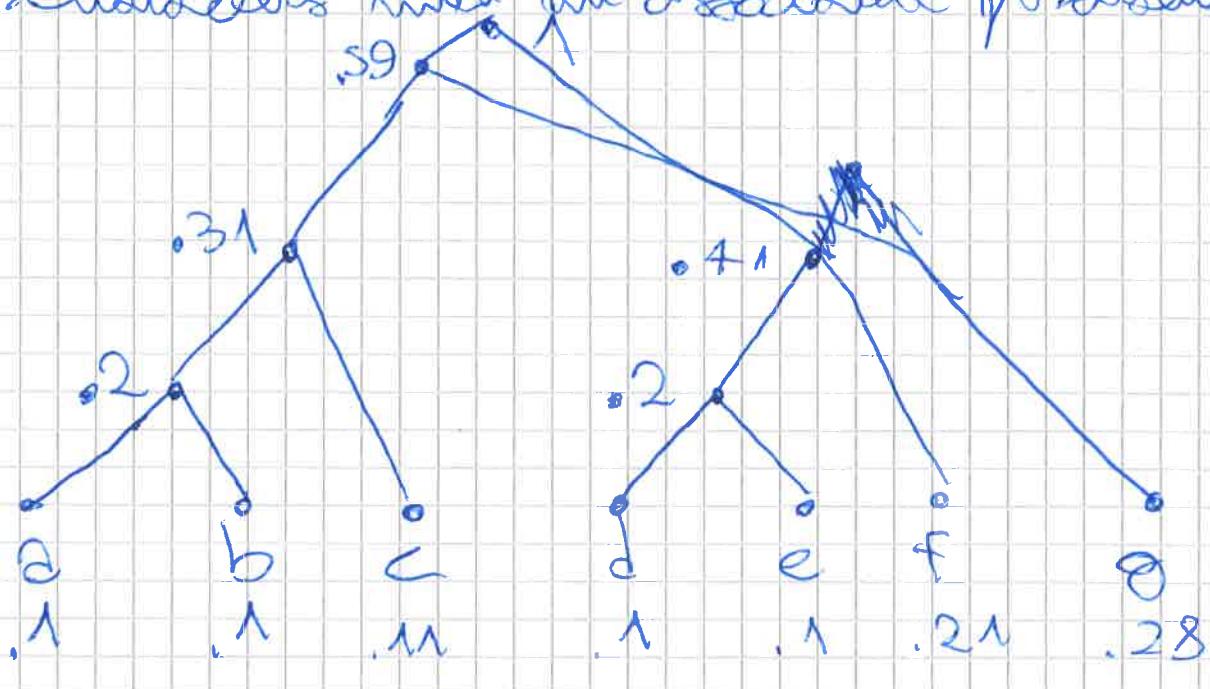
~~Idea: Keep examining~~

~~Idea: Stop the while loop when the
constraint no longer satisfies~~

~~Return $\text{SYMB}[l, V - \text{Fcv}[l]]$~~

CANONICAL HUFFMAN - COMPRESSION & DECOMPRESSION:

1) COMPRESSION: We are given the following characters with an associated probability.



① $V=0 \quad | \quad V=1 \quad | \quad V=3 \rightarrow \text{return } \text{SYMB}[2, V-8 \times l]$
 $| = 1 \quad | = 1 \quad | = 2$

$\text{SYMB}[2, 2]$

② $V=1 \Rightarrow | \quad V=2 \rightarrow \text{return } \underbrace{\text{SYMB}[2, 1]}_8$
 $| = 1 \quad | = 2$

And the symbol-table (which symbols are at which level?)

Symbol		1	2	3
	1			
1				
2	8	8		
3	c	d	e	
4	a	b		

Before decomposing, we "break" the Huffman tree into the NCL.

sym.	1	1	
0	1	2	3
2	2	1	0

$\ell = 0011011011011100$

$$\textcircled{1} \quad V=0 \Rightarrow V=0+0 \Rightarrow V=1 \quad \text{LEAVE WHILE}$$

$$I=1 \quad I=2 \quad I=3 \quad \text{return:}$$

SYMB[3, 1] = c

SYMB[3, 1] = c

$$\textcircled{2} \quad V=1 \Rightarrow V=2 \quad \text{LEAVE WHILE}$$

$$I=1 \quad I=2 \quad \text{return:}$$

SYMB[2, 2] = f

[2, 1] = 8

$$\textcircled{3} \quad V=1 \Rightarrow V=2 \quad \text{LEAVE WHILE}$$

$$I=1 \quad I=2 \quad \text{return } 8$$

\textcircled{4} return 8

$$\textcircled{5} \quad V=1 \quad V=3 \quad \text{LEAVE WHILE}$$

$$I=1 \quad I=2 \quad \text{return SYMB[l, V-fac[l]]}$$

SYMB[2, 2] = 8

ARITHMETIC ENCODING

Part of the bigger family of STATISTICAL ENCODING TECHNIQUES.

→ Recall that HUFFMAN CODING was based by the following entropy:

$$H \leq L_H + 1 < H + 1$$

So, Huffman was not ~~OPTIMAL~~ OPTIMAL for coding sources with a very small entropy (~~a very skewed distribution~~, as it uses at least one bit per symbol).

Huffman substitutes each input symbol with a codeword \rightarrow Avg. length of a TEXT compressed by Huffman is: $\Omega(1/T)$

Huffman can't be better than

$$\frac{1}{\log_2 |\Sigma|}$$

The BEST CASE is given by

A quantization of each character by 1 bit.

ARITHMETIC ENCODING:

No longer a PREFIX-CODER!



Compressed CDT is NOT a concatenation of CODE WORDS associated to symbols $\in \Sigma$, but:

"A bit ~~of~~ of the encoded output can represent more than one input symbol".

We obtain the following:

• Better compression, at the cost of slowing down the algorithm's lost capability to decode characters at any position.

[+ also works in a DYNAMIC MODEL, where probabilities $P\{O\}$ are updated as the input sequence S is processed]

HUFFMAN CODE : $H \leq L_H < H + 1$

ARITHMETIC CODE : $H \leq L_A < H + \frac{1}{n}$

BINARY REPRESENTATION AS

DYADIC FRACTION

INPUT: A bit stream ~~$b_1 b_2 \dots b_K$~~ with $K + \infty$

can be interpreted as a real number in the range $[0, 1]$ by prepending a "0" to it.

$$b_1 b_2 \dots b_K \rightarrow 0.b_1 b_2 b_3 \dots b_K = \sum_{i=1}^K b_i \cdot 2^{-i}$$

$$\sum_{i=1}^K b_i \cdot 2^{-i} = \frac{1}{2^K} = \text{DYADIC FRACTION!}$$

~~Now want to find x where:~~

~~EXAMPLE: 0.875~~

EXAMPLE: INPUT: 1101

$$\Rightarrow 0.(1101) = \frac{1}{2} + \frac{1}{4} + \frac{1}{16} = \frac{8+4+1}{16} = \frac{13}{16} = \frac{\sqrt{K}}{2^4}$$

$1 \cdot 2^{-1} + 1 \cdot 2^{-2}$

CONVERSION ALGORITHM:

REAL NUMBER \Rightarrow FRACTION (possibly INFINITE #BITS)
(Used to convert result can be converted into CANONICAL (LCFFM) AND BINARY)

According to the following algorithm:

Converter(x): $|x \in \mathbb{R} \text{ in } [0, 1]$

repeat

$$x = 2 \cdot x$$

if $x < 1$ then:

$$\text{output} = \text{output} \cup 0$$

else:

$$\text{output} = \text{output} \cup 1$$

$$x = x - 1$$

end if

until accuracy desired

|| We keep going till the desired level of accuracy is reached.

DYADIC FRACTION:

A fraction of the form $\frac{V}{2^K}$, where V and K are positive integers.

\Rightarrow Real number associated to $b_1 b_2 b_3 \dots b_K$ is the dyadic fraction $\frac{\text{val}(b_1 b_2 \dots b_K)}{2^K}$

Any fraction $\frac{v}{2^k}$ can be written as $\lim_{K} v$

BINARY REPRESENTATION
of integer v as a bit string of length K .

EXAMPLE: $x = \frac{1}{3}$

\Rightarrow Apply the converter algorithm to it.

$$1) \frac{1}{3} \cdot 2 = \frac{2}{3} < 1 \Rightarrow \text{output} = 0$$

$$2) \frac{2}{3} \cdot 2 = \frac{4}{3} \geq 1 \Rightarrow \text{output} = 01$$

~~$x = \frac{4}{3} - 1 = \frac{1}{3}$~~

The procedure could continue forever!

~~$0.0101010101 = 0.\overline{01}$~~

~~$0.33333\ldots$~~

~~1.5000000000000002~~

~~Compression by ARITHMETIC CODING~~

~~Arithmetic coding operates in intervals. It symbols~~

~~- Each step takes a subinterval $[0, 1]$ as input, representing the prefix of the input sequence compressed so far, and the PROBABILITIES of alphabet symbols Σ . The interval is further subdivided into smaller subintervals (one for each symbol σ in Σ)~~

~~for each symbol σ in Σ~~

TRUNCATION ERROR:

As we do not have an infinite # Bits at our disposal, we generally truncate somewhere our BINARY STAIR.

Suppose we truncate after the d_L -bit.
What is the error we would get?

WORST-Case ERROR:

All digits following the d_L -digit are 1.

$0.b_1 b_2 \dots b_{d_L} \overline{b_{d_L+1} b_{d_L+2} \dots}$

Truncation Error

$$x = \underbrace{0.b_1 b_2 \dots b_{d_L}}_{d \text{ DIGITS}} \overline{b_{d_L+1} b_{d_L+2} \dots}$$

\uparrow Truncation

\Rightarrow Truncating at b_{d_L} , we will get that.

$$x \in [x - 2^{-d}, x]$$

\uparrow Truncation of x

$$x = \overline{0.b_1 b_2 \dots b_{d_L} b_{d_L+1} b_{d_L+2} \dots b_{d_L+n}}$$

\uparrow difference
 x on the left
following the truncation

$$x - x = \sum_{i=1}^{\infty} 0.b_{d_L+1} b_{d_L+2} \dots b_{d_L+i}$$

$$0.b_1 b_2 \dots b_{d_L} \overline{b_{d_L+1} b_{d_L+2} \dots b_{d_L+n}}$$

$$\# \Delta_S = \sum_{i=1}^{+\infty} 1 \cdot 2^{-(d+1)} = \sum_{i=1}^{+\infty} 2^{-d-i} \quad \text{Worst-case!}$$

(Following)

$$= 2^{-d} \cdot \sum_{k=1}^{\infty} 2^{-k} = 2^{-d} \cdot \left(\sum_{k=1}^{\infty} \frac{1}{2^k} \right) = 2^{-d}$$

\Rightarrow The error reduces exponentially with the # bits used.

COMPRESSION (CODING) BY ARITHMETIC ENCODING

INPUT: - A text T :

- Alphabet Σ , with probabilities $P\{\theta\}$ associated to each character θ

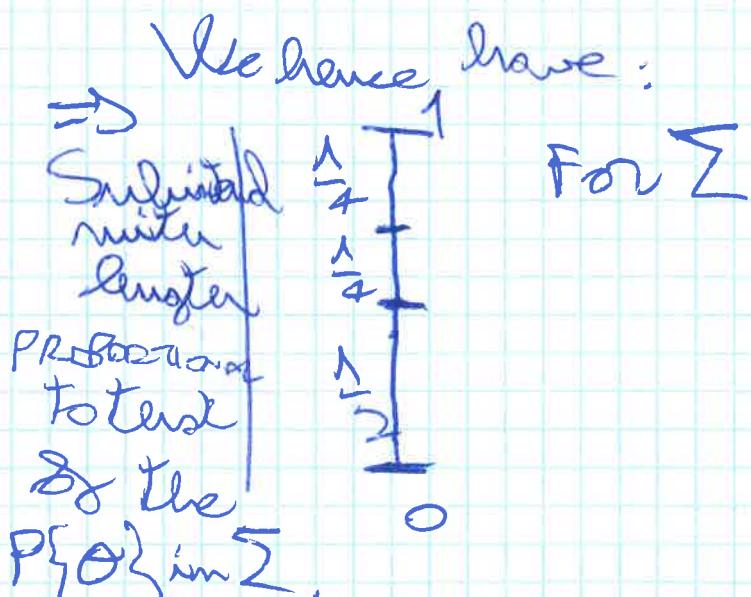
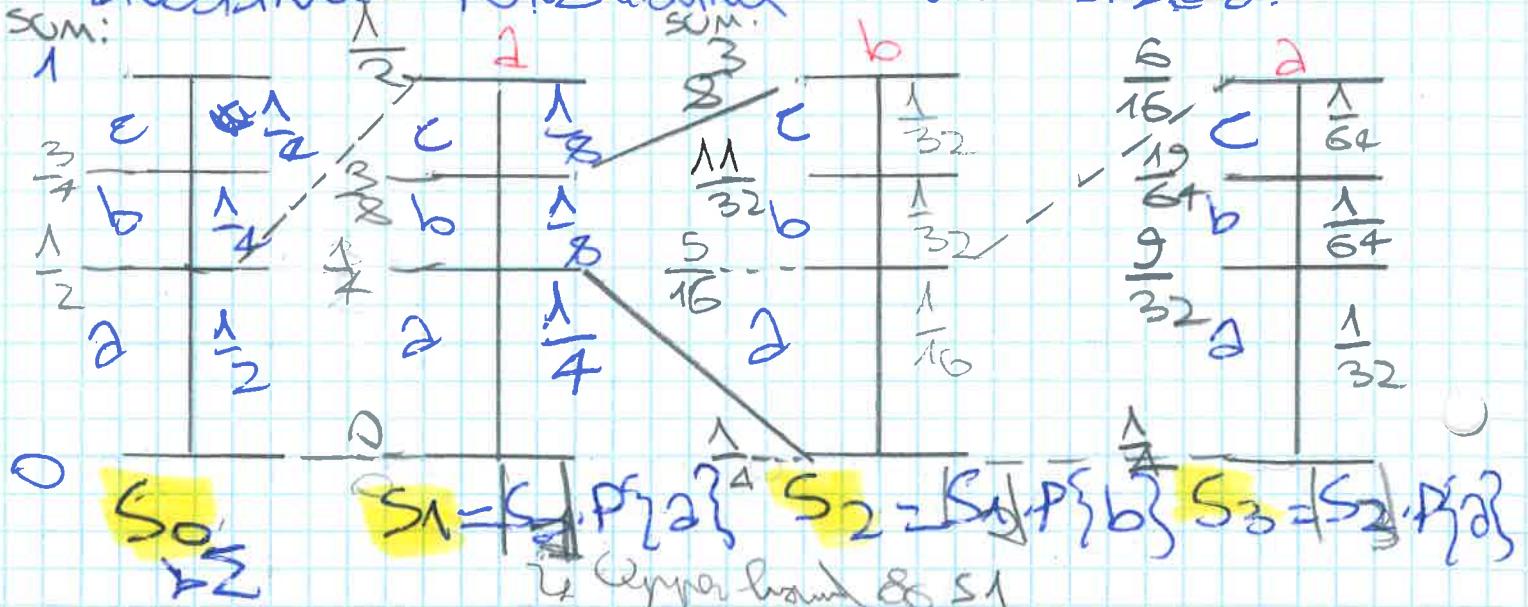
$$T = abac$$

$$\Sigma = \begin{matrix} a & b & c \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{matrix}$$

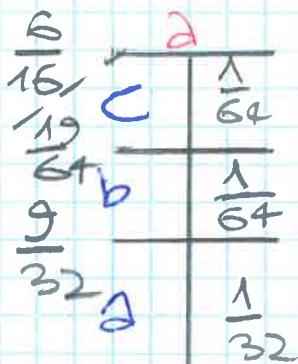
$$\Sigma = S_0$$

$$T = \begin{matrix} a & b & a & c \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} & \end{matrix}$$

ENCODING FUNCTION



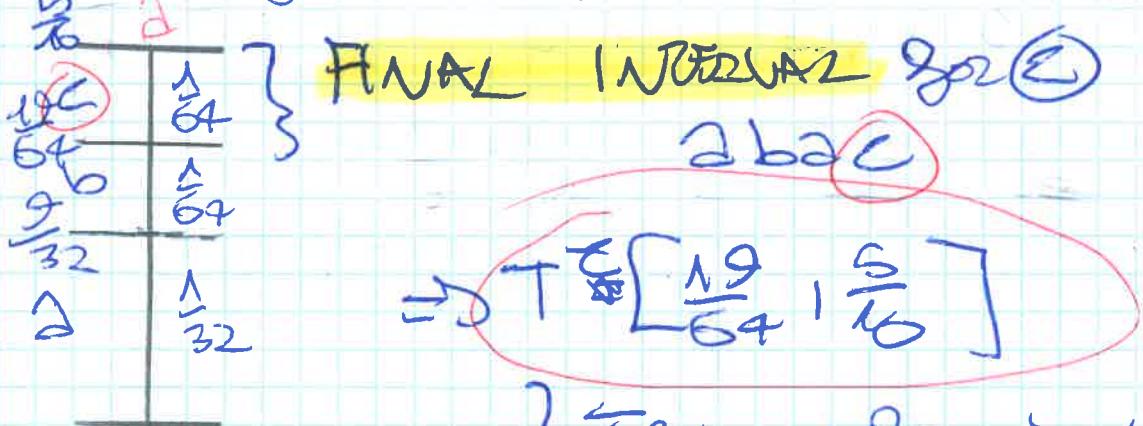
- VISUALIZED:



Upper bound of S_1

~~FIFO & Shortest Prefix~~

Process each character of the text T ,
 by expanding the corresponding Θ in
 S^T . \Rightarrow The intervals are given by these
 sum of the P_{Θ} up to text character



$2ba\cancel{c}$

$$\Rightarrow T \in \left[\frac{19}{64}, \frac{5}{16} \right]$$

RESULT

\hookrightarrow This number is inside all
 previously-generated intervals.

By

Starting from this number, together with
 (W) , we can proceed backwards and
 re-construct the ranges.

\hookrightarrow After performing 4 steps backwards,
 we stop.

DECOMPRESSION by INVERSE DECODE

INPUT:
 • Result of compression stage $\rightarrow \left[\frac{19}{64}, \frac{5}{16} \right]$
 • $|T|$

• Symbol probabilities $P\{\Theta\}_{\Sigma}$.

OUTPUT: Original sequence T .

FUNCTIONING:

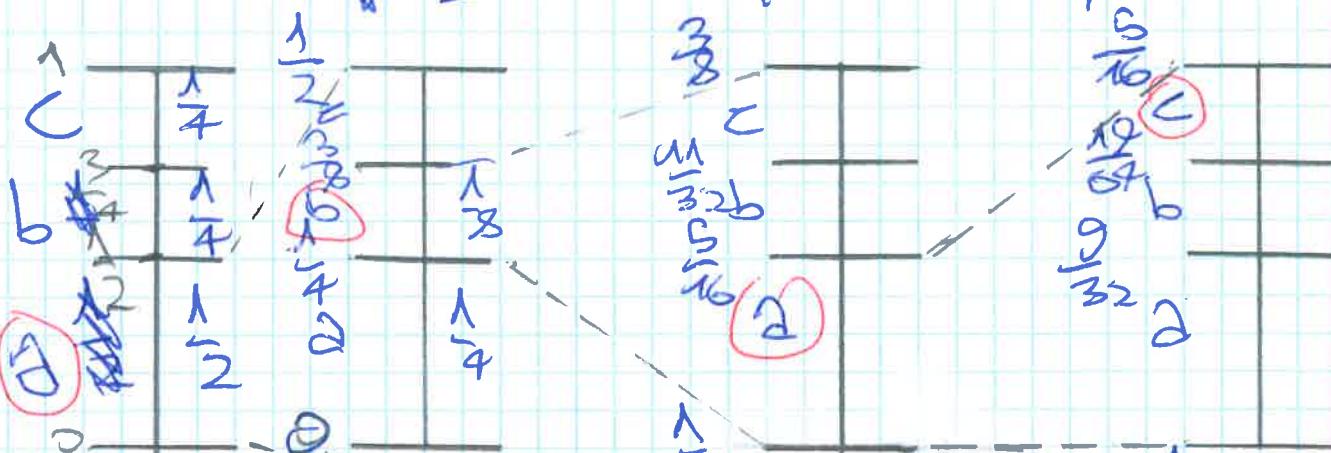
- In the first iteration, we make use of the distribution $P\{\Theta\}_{\Sigma}$.

In every other iteration, we decide one symbol using the current model; we then update the model.

EXAMPLE DECOMPRESSION.

$$\text{ENCODED RESULT} = \left[\frac{39}{128}, 14 \right] \quad \Sigma = \{a, b, c\}$$

$$P\{a\} = \cancel{\frac{1}{2}}, \quad P\{b\} = \frac{1}{4}, \quad P\{c\} = \frac{1}{4}$$



$$\text{So for } a: S_1 = S_0 \cdot P\{a\} \quad \downarrow \quad S_2 = S_1 \cdot P\{b\} \quad \downarrow \quad \frac{1}{4} S_3 = S_2 \cdot P\{a\} \quad \downarrow$$

Where does $\frac{39}{128}$ lie?
 $\frac{39}{128} \in [0, \frac{1}{2}]$
 where:
 a

Where does $\frac{39}{128}$ lie?
 $\frac{39}{128} \in [\frac{1}{4}, \frac{3}{4}]$
 where:
 b

Where does $\frac{39}{128}$ lie?
 $\frac{39}{128} \in [\frac{1}{4}, \frac{5}{16}]$
 where:
 a

Where does $\frac{39}{128}$ lie?
 $\frac{39}{128} \in [\frac{19}{64}, \frac{5}{16}]$
 where:
 c

→ The reconstructed sequence T is
 hence $T = \{a, b, a, c\}$

[After 4 stages] / we are done.
 $n=|T|$

DECOMPRESSION By ARITHMETIC ENCODING - ExFUSE.

Mixture of DYADIC FRACTIONS, not reals!

$$P\{a\} = \frac{1}{2} \quad P\{b\} = P\{c\} = \frac{1}{4}$$

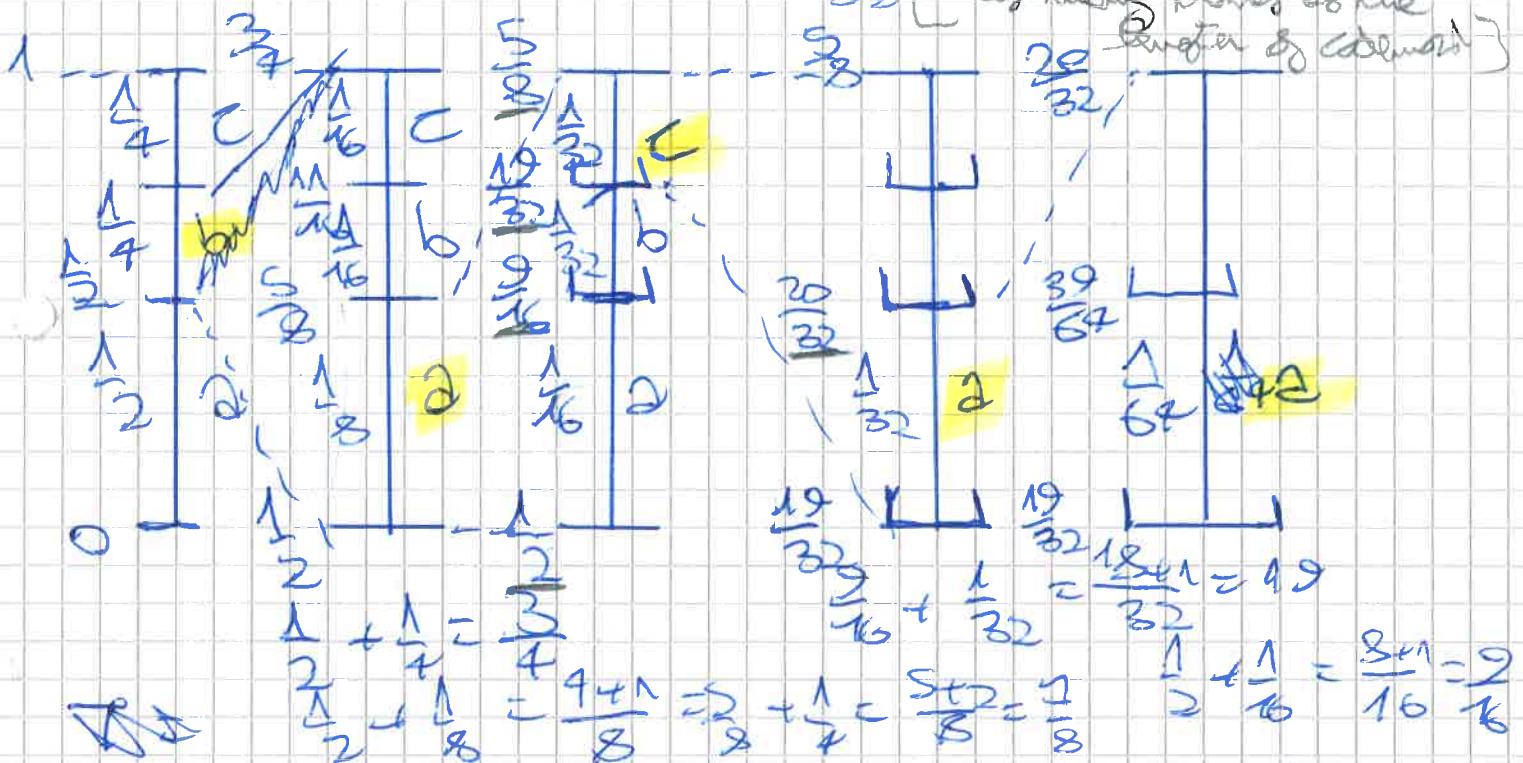
$$C = \frac{1}{2} \frac{1}{4} \frac{1}{8} \frac{1}{16} = \frac{16+2+1}{32} = \frac{19}{32}$$

First trying \rightarrow Convert into a **DYADIC FRACTION**

1) Given, model analogously to arithmetic encoding (but approximated)

2) Find out where $\frac{19}{32}$ lies freely time

[as many times as the length of codeword]



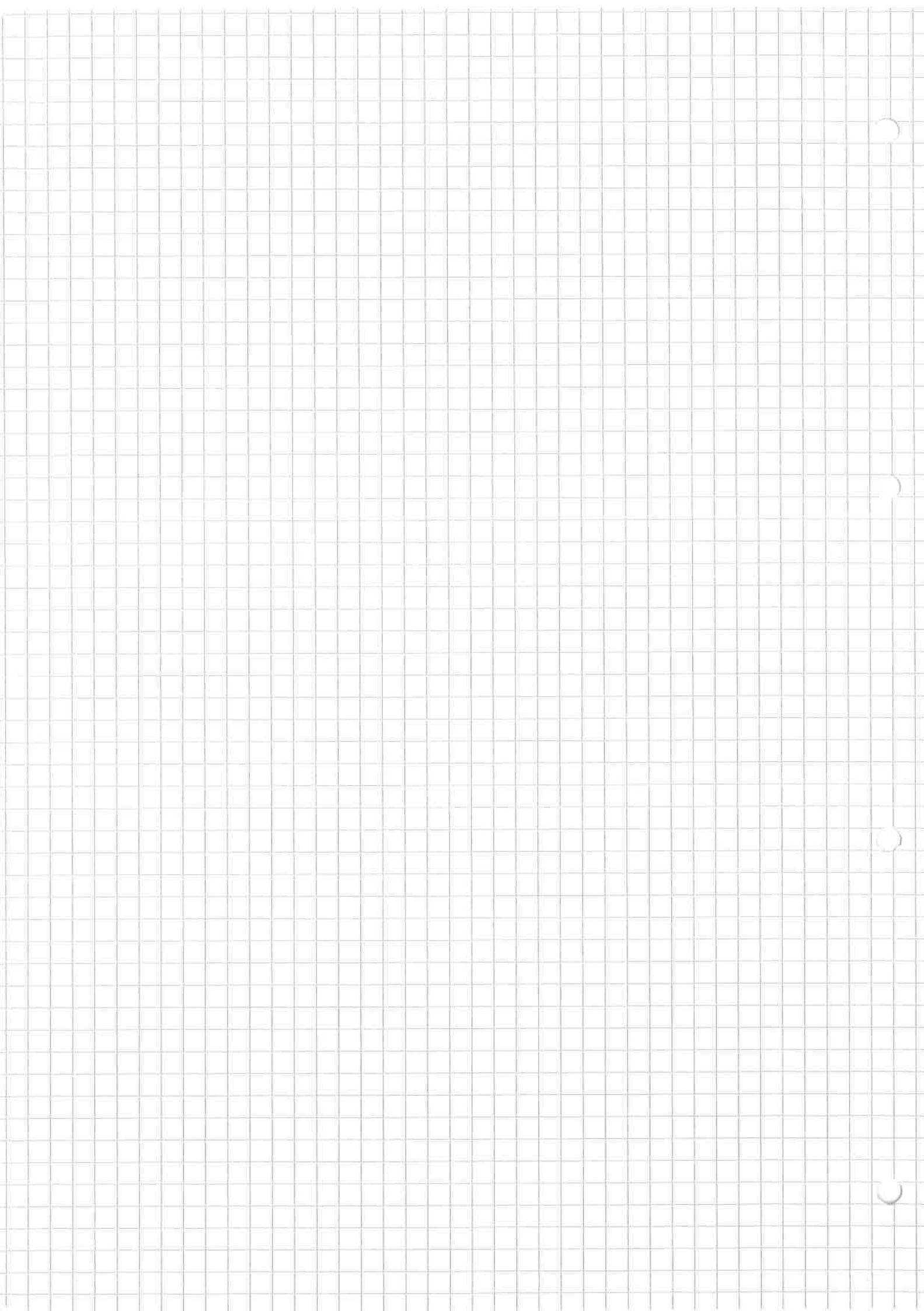
$$\boxed{T = b a c a a}$$

$$\boxed{T = b a c \overleftarrow{a}}$$

$$\frac{19}{32} + \frac{1}{64}$$

$$\frac{19}{32} + \frac{1}{32}$$

$$= \frac{38+1}{64} = \frac{39}{64}$$



EFFICIENCY of ARITHMETIC Encoding

After finding $[l_n, s_n]$ for an input sequence S , we still need to find an efficient way to encode it.

2) Encoded of $[l_n, s_n]$

(i.e.: the range from $\underline{l_n}$ to $\overline{s_n}$)
arithmetic encoding

→ We would like to pick a number in the interval $[l_n, s_n]$ such that it can be encoded with as few bits as possible.

Notice that → We have shown that:

$$X = \overbrace{X}^{\uparrow} \Sigma 2^{-d} \Rightarrow X - 2^{-d} \Sigma \overbrace{X}^{\uparrow}$$

$$\text{and } \overbrace{X}^{\uparrow} \Sigma X$$

TRUNCATION
of X after d bits

$$\text{by } s_n$$

$[l_n, \cancel{s_n}]$

Dropping bits has been set to 0

How do we encode it?

$$b = b_d \underbrace{b_{d+1} \dots b_{d+n}}_{000000} \dots$$

ENCODING By "TRUNCATION"

$[l_n, s_n] \Rightarrow [l_n, l_n + s_n]$

~~Truncating $l_n + s_n$ to its~~

$\log_2 \frac{2^n}{s_n}$ bits falls in $[l_n, l_n + s_n]$

Goal: $l = l_n; s = s_n$

$[l, s] \Rightarrow [l, l+s]$ can be truncated. Where?

\Rightarrow Truncate $l + \frac{s}{2}$ to its first $\lceil \log_2 \frac{2}{s} \rceil$

This would make it "full" in the interval $[l, l+s]$

PROOF: Set $d = \lceil \log_2 \frac{2}{s} \rceil$

$$2^{-d} = \frac{1}{2^d} \leq \frac{1}{2^{\lceil \log_2 \frac{2}{s} \rceil}} \leq \frac{1}{2} \leq \frac{s}{2}$$

$$\Rightarrow 2^{-d} \leq \frac{s}{2}$$

RETURN VALUE: $s_{1:2^d}$

\Rightarrow ARITHMETIC Encoding to hence
goes to return the first $\lceil \log_2 \frac{2}{s} \rceil$ bits
(TRUNCATION)

of the BINARY REPRESENTATION of value:

$l_n + \frac{s_n}{2}$ with $\lceil \log_2 \frac{2}{s} \rceil$ bits.

EXAMPLE: We ~~had~~ found test the resulting
interval was: $[\frac{19}{64}; \frac{5}{16}]$

$$[\frac{19}{64}; \frac{20}{64}] = [l_n, s_n]$$

$$\Rightarrow \left[\frac{19}{64} : \frac{19}{64} + \frac{1}{64} \right] \quad l_n = \frac{19}{64} \\ l_n \quad l_n \quad s_n \quad s_n = \frac{1}{64}$$

\Rightarrow OUTPUT decomposition shows

Binary representation of $l_n + \frac{s_n}{2}$

$$l_n + \frac{s_n}{2} = \frac{19}{64} + \frac{1}{64} \cdot \frac{1}{2} \\ = \frac{19}{64} + \frac{1}{128} = \frac{38+1}{128} = \frac{39}{128}$$

\Rightarrow We now need to represent $\frac{39}{128}$ by making use of $\log_2 \frac{2}{5}$ bits.

$$\log_2 \frac{2}{\frac{5}{128}} = \log_2 \frac{128}{5} = 7$$

\Rightarrow We will make use of 7 bits!

How to find these 7 bits?

\Rightarrow Convert $\frac{39}{128}$ in binary (Converter Algorithm)
and truncate at the 7th bit!
7 times!

$$\Rightarrow \frac{39}{128} \cdot 2 < 1 \Rightarrow \text{output} = 0$$

$$\frac{39}{128} \cdot 2 < 1 \Rightarrow \text{output} = 1$$

$$\left(\frac{39}{128} - 1 \right) \cdot 2 < 1$$

$$\frac{56}{128} < 1 \Rightarrow \text{DOUBLER} = 0$$

$$\frac{112}{128} \cdot 2 < N^2 \Rightarrow \text{DOUBLER} = 1$$

$$\left(\frac{224}{128} - 1\right) \cdot 2 = \frac{192}{128} \geq 1 \Rightarrow \text{DOUBLER} = 1$$

$$\left(\frac{192}{128} - 1\right) \cdot 2 = 1 \geq 1 \Rightarrow \text{DOUBLER} = 1$$

\Rightarrow FINAL OUTPUT:

$$\underbrace{01001111}_{\text{Binary Rep.}} , 4 \xrightarrow{\text{SST}} 4|T|$$

PROOF of the RELATION of #BITS and ENTROPY. (i.e. it is no 2 BITS from the OPTIMAL)
 i.e. we have 2 BITS from the OPTIMAL
 LOWER BOUND

(e.g. Relation with OPTIMAL #BITS)

We know that the ~~SIZE~~ SIZE of a string $S_n = \prod_{i=1}^n P\{T[i]\}$ Line. at pos. i

And we have seen that S_n will be represented by $\log_2 \frac{2}{S_n}$ BITS.

$$\lceil \log_2 \frac{2}{S_n} \rceil < \overbrace{\log_2 \frac{2}{S_n}}^{\text{# Bits at which we truncate ARITHMETIC CODEC}} + 1$$

Bits at which we truncate ARITHMETIC CODEC

$$S_n = \prod_{i=1}^n P\{T[i]\}$$

$$2 - \log_2 S_n$$

$$2 - \log_2 \left(\prod_{i=1}^n P\{T[i]\} \right)$$

By Logarithms law

$$\begin{aligned} \log_2(a \cdot b) \\ = \log_2 a + \log_2 b \end{aligned}$$

$$2 - \sum_{i=1}^n \log_2 P\{T[i]\}$$

$n\theta = \# \text{occurrences}$

$$\begin{aligned} P\{\theta\} = \frac{n\theta}{n} & \quad \text{The group } \theta \text{ - some symbol} \\ & \quad \text{is in } n\theta \\ & \quad * \text{ suff. long} \end{aligned}$$

$$2 - \sum_{\theta}^{n\theta} \log_2 P\{\theta\}$$

Div & mult. by n .

$$2 - n \cdot \sum_{\theta} \frac{n\theta}{n} \cdot \log_2 P\{\theta\}$$

Faktor n !

$$2 - n \cdot \sum_{\theta} P\{\theta\} \cdot \log_2 (P\{\theta\})$$

We know the DEFINITION of ENTHROPY:

$$H = \sum_{\theta} P\{\theta\} \cdot \log_2 \left(\frac{1}{P\{\theta\}} \right)$$

ARITHMETIC ENCODING

$$\cancel{- \log_2 \frac{2}{S_n} \Rightarrow \lceil \log_2 \frac{2}{S_n} \rceil < 2 + n \cdot H_0}$$

Now, we can get the length of every symbol by dividing by n .

$$\text{Length } L_A = \frac{nH}{n} + \frac{2}{n}$$

$$L_A = H + \frac{2}{n}$$

Or:

$$\Rightarrow \lceil \log_2 \frac{2}{S_n} \rceil < 2 + n \cdot H$$

\Rightarrow There is a waste of just 2 bits for an input sequence S , or $\frac{2}{n}$ bits per symbol of the sequence!

These 2 bits are hence negligible for longer and longer sequences.

\Rightarrow HUFFMAN:

$$n + n \cdot H$$

ARITHMETIC:

$$2 + n \cdot H$$

CANONICAL is faster than A. Better than H.
& can decompress any portion.

+ On-the-fly!

Used when not using streams
DISCRETE

Only entire file is compressible.

Generally used when highly skewed DISR.

(13) DICTIONARY-BASED COMPRESSORS

Different approach from the statistical one! \Rightarrow Not interested in PROBABILITIES of source S!

~~S = "Cane"~~ D = "cane" \rightarrow "z" } Dictionary & MAPPING
"gatto" \rightarrow "o" } (toestring replaced by a
"mento" \rightarrow "za")

S = "Cane", "mento", "gatto" via a
S' = "z", "za", "o" Italian

The DICTIONARY is of GREAT IMPORTANCE for achieving good compression!

(Ex: Italian text dictionary w/ English text)

ZIV & ZEMPEL (Z - xx) introduces a series of compressors, where each string occurrence is substituted with:

- Either the offset of the previous occurrence
- ID assigned incrementally to new dictionary phrases \Rightarrow DYNAMIC DICTIONARY.

Starts empty and grows
(beginning of input sentence to process it)
slow compression, but after some steps good compression is obtained.

A 203% compression ratio.

L277:

Performance comparable to that of ARITHMETIC & MUFFMAN coding (according to the author).

ALGORITHMIC GROUND BASE:

Sliding Window $W[1, W]$, containing a portion of the input sequence that has been processed so far [typically last W characters] and look-ahead buffer B , which contains suffix of the text till to be processed.

Ex: ~~babababaaabb~~

$W = aabbbaabb$

$B = baababaaabb$

$\underbrace{aab}_{W} \underbrace{aabab}_{B}$

Algorithm processes in 2 phases:

— PARsing: Transform the input sequence into a sequence of triples of integers.

— ENCoding: Turns triples into a compressed bit-stream by applying a statistical compressor and integer encoder scheme.

FUNCTIONS VISUALIZED.

$\text{V} | \text{aabbabab} \Rightarrow \langle 0, 0, 2 \rangle$

$\text{V} | \text{aabbabab}$ distance of a from the start = length
copy $\Rightarrow \langle 1, 1, b \rangle$

$\text{V} | \text{aabbabab}$ $\Rightarrow \langle 1, 1, 2 \rangle$

$\text{V} | \text{aabbabab}$ $\Rightarrow \langle 2, 3, \text{End of string} \rangle$

PARTIAL STATE'S FUNCTIONALITY:

Search for the longest prefix d of B that occurs as a substring of VB (not B , as the previous occurrence ^{CONCENATION (either VB or B)} we are searching for starts in VB and extends up to B).

\Rightarrow If d occurs at distance d from cur. position (beginning of B) and it is followed by c , then we emit:

$\langle d, |d|, c \rangle$

distance of d from the beginning of B . $|d|$ char following d

$\Rightarrow \text{VB} | \text{aabbabab} \Rightarrow \langle 0, 0, B[1] \rangle$ is not match!

$\text{V} | \text{aabbabab} \Rightarrow \langle 0, 0, 2 \rangle$

$\text{V} | \text{aabbabab} \Rightarrow \langle 1, 1, b \rangle$

\Rightarrow If a match of d units of B in VB is found, we extend advance the window by $|d| + 1$ positions in B and slides VB correspondingly.

\Rightarrow If a match is NOT found, the output triple is $\langle 0, 0, B[1] \rangle$.

~~SLIDING WINDOW~~: Delimits the size of the dictionary (quadratic in W 's length)

~~W=4~~ $\underline{aabbabab} \Rightarrow \langle d, 1d \rangle$
 $\Rightarrow \langle 0, 0, 2 \rangle$

$\{$ $\underline{a|abbabab} \Rightarrow \langle 1, 1, b \rangle$
 $\quad W=2$ $\quad \underline{b}$

$\underline{aabb|babab} \Rightarrow \langle 1, 1, 2 \rangle$

$\underline{aabbba|bab} \Rightarrow \langle 2, 3, Ed \rangle$

4 Performance of LZ77 is strongly depend on W 's length [exposure]

- + The shorter W , the worse the compression ratio [longer compression time]
- + The larger W , the better the compression ratio [but the longer compression time]

(1982) LZSS: In this version, we improve over LZ77, by wanting to add the short tail (following if it's not found) we ~~do~~ not limit two 0s all the time.

⇒ Solution: Just limit:

• If a ~~0~~ $\in WB \subset \langle d, 1d \rangle \Rightarrow$ Advanced by 1d

• B a ~~not~~ $\in WB \subset \langle 0, B[1] \rangle$

EXAMPLE :

B distance of d from the beginning of B.

1 2 2 bbaabab $\Rightarrow \langle 0, 2 \rangle$

1 2 1 abbababba $\Rightarrow \langle 1, 1 \rangle$

2 2 1 bbaabab $\Rightarrow \langle 0, b \rangle$

2 2 1 bbaabab $\Rightarrow \langle 1, 1 \rangle$

2 2 bbbabb $\Rightarrow \langle 3, 2 \rangle$

2 2 bbabab $\Rightarrow \langle 2, 2 \rangle$

\Rightarrow We can then encode the emitted pairs by STATISTICAL Encoding or by INTEGERS Encoding.

[GZIP is based on LZSS]
Fresh tables

LZ77.

Q: DISTANCE & compute

Sliding window of LZ77 stores up to the longest phrase to encode, but also limits the storage space (and hence the ultimate compression ratio)

\Rightarrow Build INCREMENTALLY an EXPLICIT DICTIONARY that contains a subset of strings of the input sequence S_1 .

$D =$ Dictionary where every phrase s_j is identified via integer ID(j).

\Rightarrow We then probe S_1 and determine the longest prefix s_i^* that is also a phrase of D .

→ We then substitute the s -grate phrase by $\text{pair} \leftarrow \text{ID}(f')$, where f' is the character following f in s .

→ String pairs will then be encoded by STATISTICAL OR VARIABLE LENGTH INTEGERS encoding.

EXAMPLE D:

Phrase	ID:
abcd	43
bcde	44

$S = \underline{ab}cd\text{efg}$

- ① Output: $\leftarrow 43, e \rightarrow$
 \Rightarrow Adds $ab\text{cd}\text{e}$, $\rightarrow 44$ to D.

EXAMPLE 2:

- We construct a TIE based on the BIGRAMS output.

EXAMPLE 2. S' = Sequence yet to be parsed.

D = Current dictionary, where every phrase δ is identified via integer $i\delta/g$.

Parse δ of S' (part ~~not~~ yet to be parsed). Consists by determining the longest prefix δ' that is also a phrase of D , and substituting it with the pair: LEAF: $\langle i\delta/\delta' \rangle, \langle \cdot \rangle$, where $\langle \cdot \rangle$ is the character following δ' in S' .

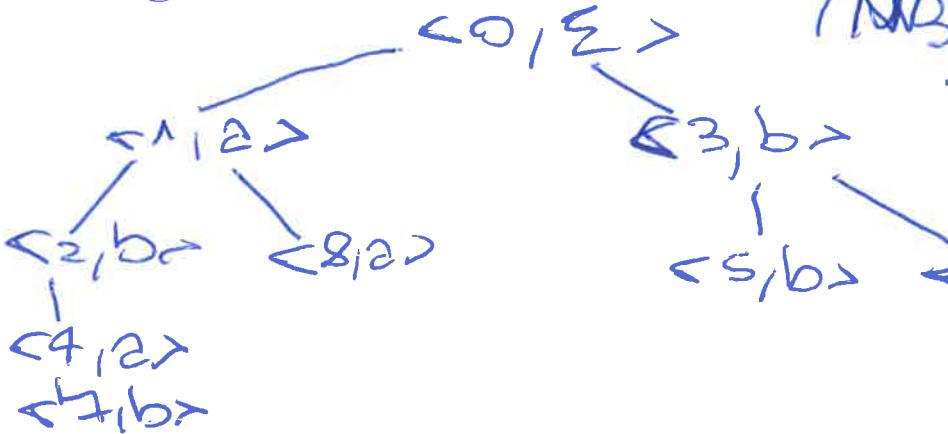
Looked up in the dictionary

PHRASE:

$S = \text{gabbababbabababaa}$ $\leftarrow \delta^{\text{incomplete}}$

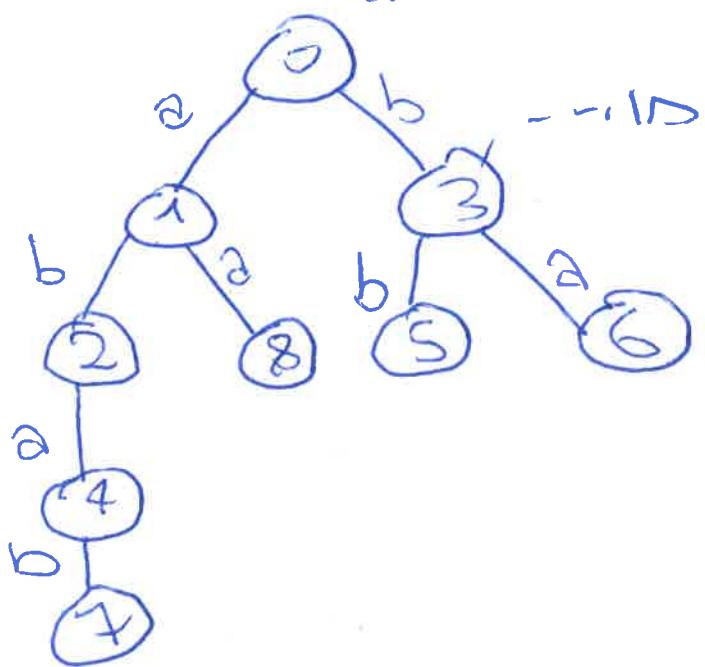
PHRASE:	INDEX	PREFIX	DICTIONARY	LEAF TO ADD
String δ contained in the dictionary D .	\emptyset	$\langle \emptyset, \emptyset \rangle$	$0: \emptyset$	$\langle 0, \emptyset \rangle \vee$
	a	$\langle 0, a \rangle$	$1: a$	$\langle 1, a \rangle$
	ab	$\langle 1, b \rangle$	$2: ab$	$\langle 2, b \rangle$
	b	$\langle 0, b \rangle$	$3: b$	$\langle 3, b \rangle$ $\langle 1, b \rangle$
	aba	$\langle 2, a \rangle$	$4: aba$	$\langle 4, a \rangle$
	bba	$\langle 3, b \rangle$	$5: bba$	$\langle 5, b \rangle$
	ba	$\langle 3, a \rangle$	$6: ba$	$\langle 6, a \rangle$
	$abab$	$\langle 4, b \rangle$	$7: abab$	$\langle 7, b \rangle$
	aa	$\langle 1, a \rangle$	$8: aa$	$\langle 8, a \rangle$

→ Based on this table, we can now build the corresponding TRIE.



NB: We construct the Trie based on the occurrence of the leaf nodes ($\langle \cdot \rangle$) in the DICT.

This tree can also be represented in an UNCOMPRESSED FASHION:



GZIP - SMART & FAST IMPLEMENTATION OF LZ77 [HASH-TABLE-based]

Key Problem: Fast search needed for longest prefix of B that repeats in V .

→ "BRUTE-FORCE APPROACH": Check the occurrence of every prefix of B in V via a linear BACKWARD SCAN, every single time.

→ HASH-TABLE-based APPROACH.

We ~~can't~~ make use of a HASH TABLE to determine α and find its previous occurrence in V .

Key Idea: Store all 3-grams occurring in V / the sliding window, namely all triplets of contiguous characters.

~~S = ABABBCABA
ABAABBAA
1 2 3 4 5 6 7 8 9 10~~

EXAMPLE:

~~S = ababbcabab ...~~

HASH TABLE:

3-GRAM	POSITION IN THE STRING
aba	0 + 5
bab	1 + 6
abc	2 + 7
bca	3 + 8
cab	4 + 9

~~POSITION OF THE 3-GRAM IN THE STRING~~

~~Position of the 3-gram in the string~~

The 3-GRAMS act as "ANCHORS" to be searched for in the table.

→ Hash table stores all the multiple ^{positions} ~~occurrences~~
renders of the 3-Gram at its corresponding
key [Searched by increasing ~~values~~
positions]

→ When X shifts to the right ~~now~~,
we output $\langle \text{d}, l \rangle$, and we delete the d
^{3-Gram shift to}
3-Grams ~~positions~~ at $\text{X}^{[1, l]}$, and insert
the new 3-Grams ~~positions~~ at $\text{B}^{[1, l]}$.
(i.e. delete the 3-Grams found so far, up to
 d , and insert new 3-Grams after d).

Using 3-Grams?

- ↳ If 2-(bit) -Gram: We would have to consider too many positions
- ↳ 4-Gram: We would have lots of overlaps.

Search (for a 3-Gram a) in X .

↳ First, search $\text{B}^{[1, 3]}$ in the hash table.

} ↳ If $\text{B}^{[1, 3]}$ does not occur in the hash table, then (2P units) $\langle \text{d}, \text{B}^{[1]} \rangle$ and padding advances by one char.

↳ If $\text{B}^{[1, 3]}$ does occur in the hash table, then the list L of occurrences of $\text{B}^{[1, 3]}$ in X is determined.

↳ Once an L is found (i.e. array of positions of $\text{B}^{[1, 3]}$ in S), for each position (i) in L the alg. compares character-by-character

$S[i, n]$ compares with B to compute the LCP (Longest Common Prefix).

→ Found i^* $\in L$ (i.e. list of characters sharing longest common prefix)
 $i^* = \text{INDEX}$ of α

• P = current position of B in S .

→ Algorithm finds pair:

$(P - i^*, |\alpha|)$

[Goal: Find the longest α]

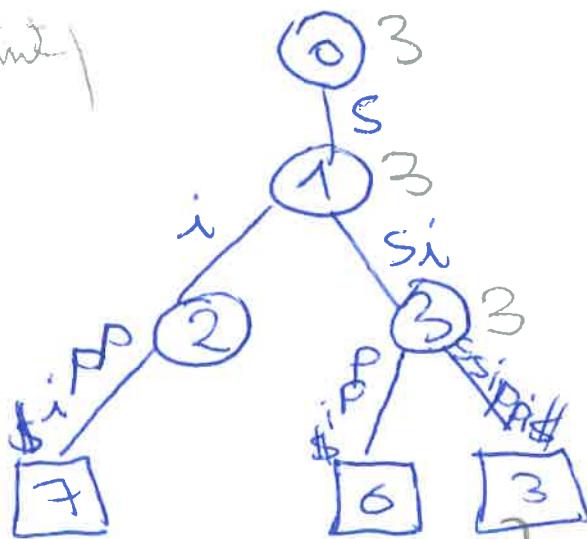
OPTIMAL SOLUTION: Based on the SUFFIX TREE of the entire text (actually not used).
→ Balance between theory ^{in practice} & practice.

[Trade-off ^{ENGINEERING COMPROMISES} between optimality & practicality]

SUFFIX TREE-based APPROACH:

① Build suffix tree over T . $T = \text{mississippiy\$}$

(only relevant part)



LONGEST SUFFIX that repeats

(2) Pre-compute the MINIMUM in every single leaf ~~to find~~ ^{levels} among a node's children.

When searching for a value in the SUFFIX TREE, we keep going down as long as $\text{minimum} < n$. [i.e.: find leaf node] (i)

TIME: $O(n)$ → Go down ~~root node~~ as much as string's length.

SPACE: $O(n)$ → Needs to store all suffixes of T.

#COMPARISONS proportional to a string's length.

i = current position.

The limit:

~~String, char
Distance from
NNN~~ character length

~~String, |d|, ↗~~

~~↳ character following d~~

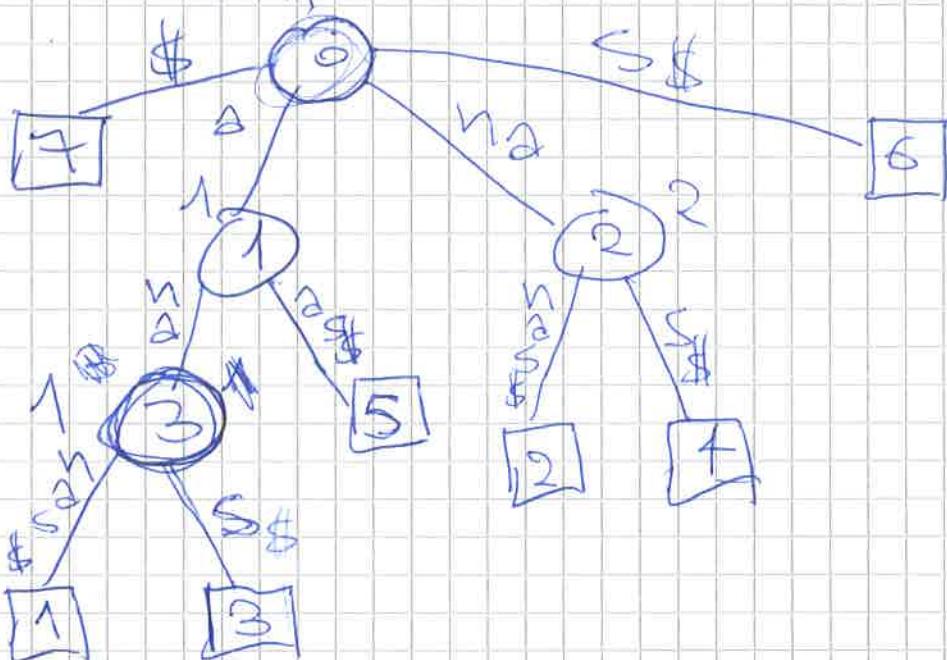
~~Distance
from the minimum
characters considered~~

→ Index of d in T.

L2 - PARSENA hash on SUFFIX TREES:

- BUILD SUFFIX TREE hash on the input string + \$

$T = \text{Paranash} \$$



- We are given a DATA STRUCTURE that allows us to compute the ~~MIN~~ LCA of two nodes among its descendants.

- We can now compute the LCA for each character (go down as long as ~~MIN~~ ~~MIN~~ ~~MIN~~).

~~MIN~~ ~~MIN~~ ~~MIN~~

$\alpha = \emptyset$ $\alpha = \alpha_1$ $\alpha = \alpha_1 \alpha_2$

\downarrow \downarrow \downarrow

$\langle P, Q, R \rangle$ $\langle P, Q, R \rangle$ $\langle P, Q, R \rangle$

distance of α from beginning \downarrow \downarrow

$\langle \text{DISTANCE} \rangle$ $\langle \text{DISTANCE} \rangle$ $\langle \text{DISTANCE} \rangle$

\downarrow \downarrow \downarrow

$\langle 1, 0, 1, 2 \rangle$ $\langle 0, 0, 1, 2 \rangle$ $\langle 3, 1, 3, 5 \rangle$

$\langle 2, 3, 5 \rangle$

$\langle \text{DISTANCE} \rangle$ $\langle \text{DISTANCE} \rangle$ $\langle \text{DISTANCE} \rangle$

\downarrow \downarrow \downarrow

$\langle 1, 1, 1, 1 \rangle$ $\langle 1, 1, 1, 1 \rangle$ $\langle 1, 1, 1, 1 \rangle$

$\langle 1, 1, 1, 1 \rangle$

N.B.: Same result as "standard" L2 FT PARSENASH \rightarrow Go down as long as $\text{MIN}(n)$ characters match with α

NORMAL L2 ITA Parsing:

$$T = \text{d} \underline{\text{h}} \text{a} \underline{\text{n}} \text{a} \underline{\text{s}} \#$$

(1) $T = | \overbrace{\text{d} \underline{\text{h}} \text{a} \underline{\text{n}} \text{a} \underline{\text{s}}}^B |$

$\mathcal{W} = \emptyset \Rightarrow$ longest prefix of \mathcal{B} is empty.
This is a suffixing of V.B.
 $\Rightarrow < \mathcal{D} | \mathcal{H} | \mathcal{A} >$

(2) $T = \text{d} \underline{\text{h}} \text{a} \underline{\text{n}} \text{a} \underline{\text{s}}$

$$\mathcal{W} = \mathcal{A}$$

$$\mathcal{L} = \emptyset$$

$$\Rightarrow < \mathcal{D} | \mathcal{H} | \mathcal{A} >$$

(3) $T = \text{d} \underline{\text{h}} \text{a} \underline{\text{n}} \text{a} \underline{\text{s}} |$

$$\mathcal{W} \quad \mathcal{B}$$

$$\mathcal{L} = \mathcal{D} \mathcal{H} \mathcal{A}$$

$$\Rightarrow < \mathcal{D} | \mathcal{H} | \mathcal{B} | \mathcal{S} >$$

We keep going down as long as,

(and rules match along the way)

$$m \leq i$$

, where:

$i =$ starting pos. of \mathcal{L} in T .

$m = \min(\mathcal{V}) \Rightarrow$ value of
the minimum at position n .

RETURN:

$$< i - m, | \mathcal{L} |, >$$

TREAPS:

14 AM

A TREAP is a BINARY SEARCH TREE with a modified way of ordering the nodes.

Each node has a ~~key~~ (BST PROPERTIES) and a priority (min-MAX PROPERTIES).

Note: $\frac{\text{Max key}(x)}{x}$

x is priority(x)

Assumption: All priorities & all keys are DISTINCT as from one another.

Hence above two nodes respect the following properties,

If V is a ~~left~~ child of U , then:

$$\boxed{\text{key}(V) < \text{key}(U)}$$
 BST Properties

If V is a right child of U , then:

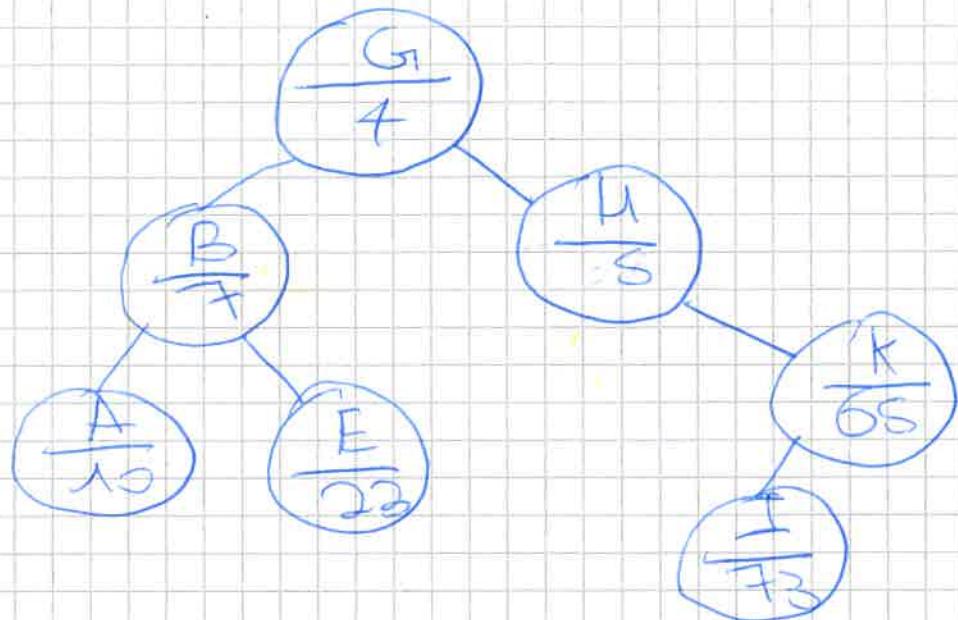
$$\boxed{\text{key}(V) > \text{key}(U)}$$
 BST Properties

If V is a child of U , then: max min
MAX MIN
priority

$$\boxed{\text{priority}(V) > \text{priority}(U)}$$

WNR for TREAPS If we insert nodes x_1, x_2, x_n with associated keys into a TREAP, then the resulting TREAP is the tree that would have been formed if nodes had been inserted into a normal BST in the order given by their $\text{priority}[x_i] < \text{priority}[x_j] \rightarrow$ node x_i inserted before x_j .

EXAMPLE TREAP:



OPERATIONS ON TREAPS

- INSET(K): Used for building the Treap, one element at a time.

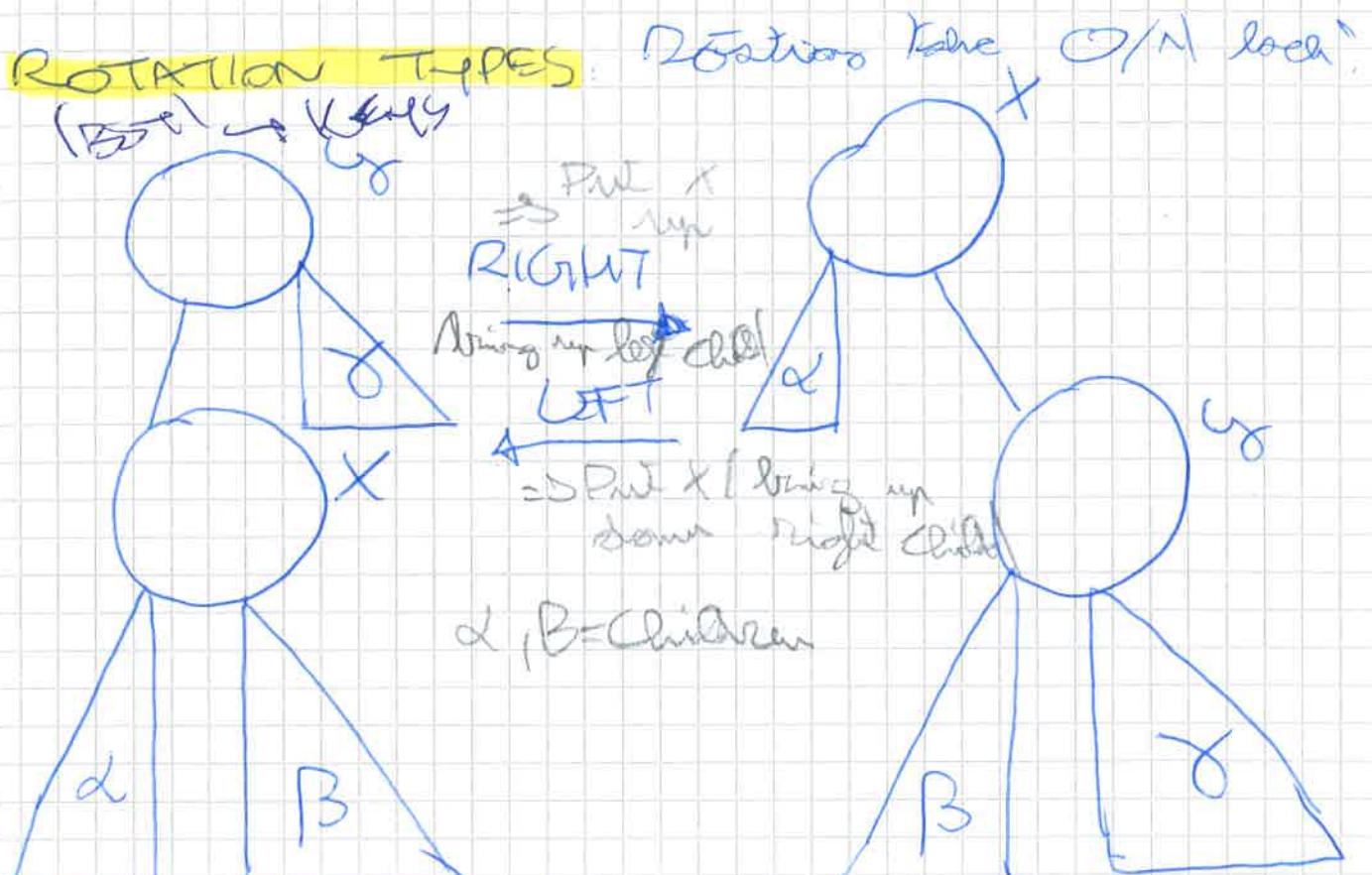
↳ Fundamental STANDARD BST INSERTION

- (1) Follow the ~~bst~~ to insert node ~~new~~ down to the leaf (at a level)
- (2) Fix the min-max property, ensuring the priority property is respected.

→ O(n)

↳ ROTATION: As long as priority (c) ≠ priority (p), we perform a rotation, which
increases/decreases parent (p)'s depth by 1.
Get left child → Decrease child (c)'s depth
Get right child → Increase child (c)'s depth by 1.

And yet, we still maintain BST's properties while doing this.



- **SEARCH(k)**: Need to navigate whole tree up to max. depth / the tree is NOT necessarily balanced.

→ Operates like ~~SEARCH~~ over a BINARY SEARCH TREE:

- [Priority] Priority $\leq \Theta$.

Bound the priority (i.e.: go down as long as the condition is satisfied).

NB! We count search by a priority being though!

Because we are dealing with a MIN-HEAP data structure.

EXAMPLE . INSERTION:

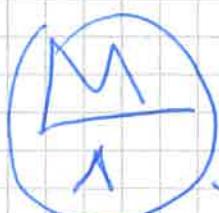
We want to insert $\frac{R}{2} = k$ into:

1) Search for

The right LEAF position

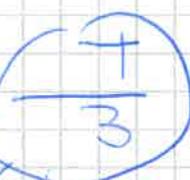
in the tree

& add k there.



= k

into:

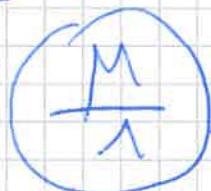


2) Need to perform

ROTATIONS to ensure

BIN SEARCH structure is preserved ; as well as
MIN. HEAP property LEFT

\rightarrow perform a ~~LEFT~~ RIGHT ROTATION.



stop



\Rightarrow perform a (RIGHT ROTATION)

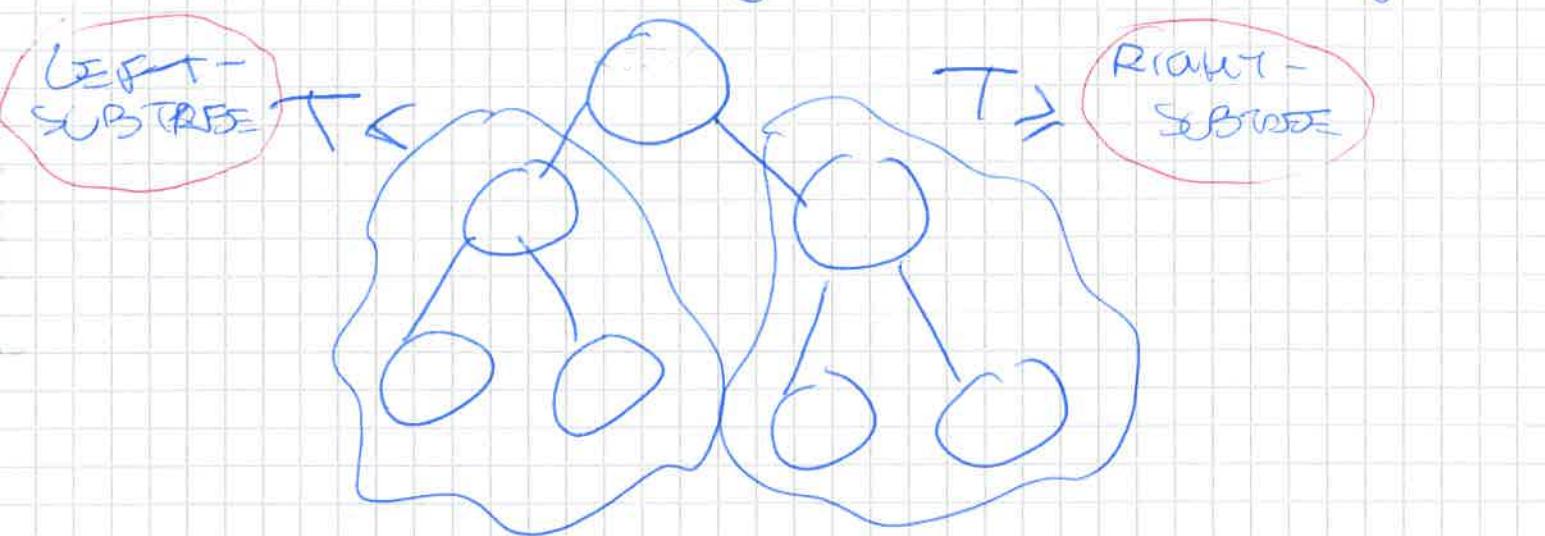


- **DELETE (K)**: Opposite of INSERT (k).

\Rightarrow Push the node down to make it a leaf via ROTATIONS \Rightarrow Then remove it safely.



- **SPLIT (K)** \Rightarrow If $K == \text{Root} \Rightarrow$ Easy!



- If $K \neq \text{root} \Rightarrow$ No Priority To (-D)

$$\Rightarrow \text{Ex: } (\underline{R}) \Rightarrow (\underline{R})_{-D}$$

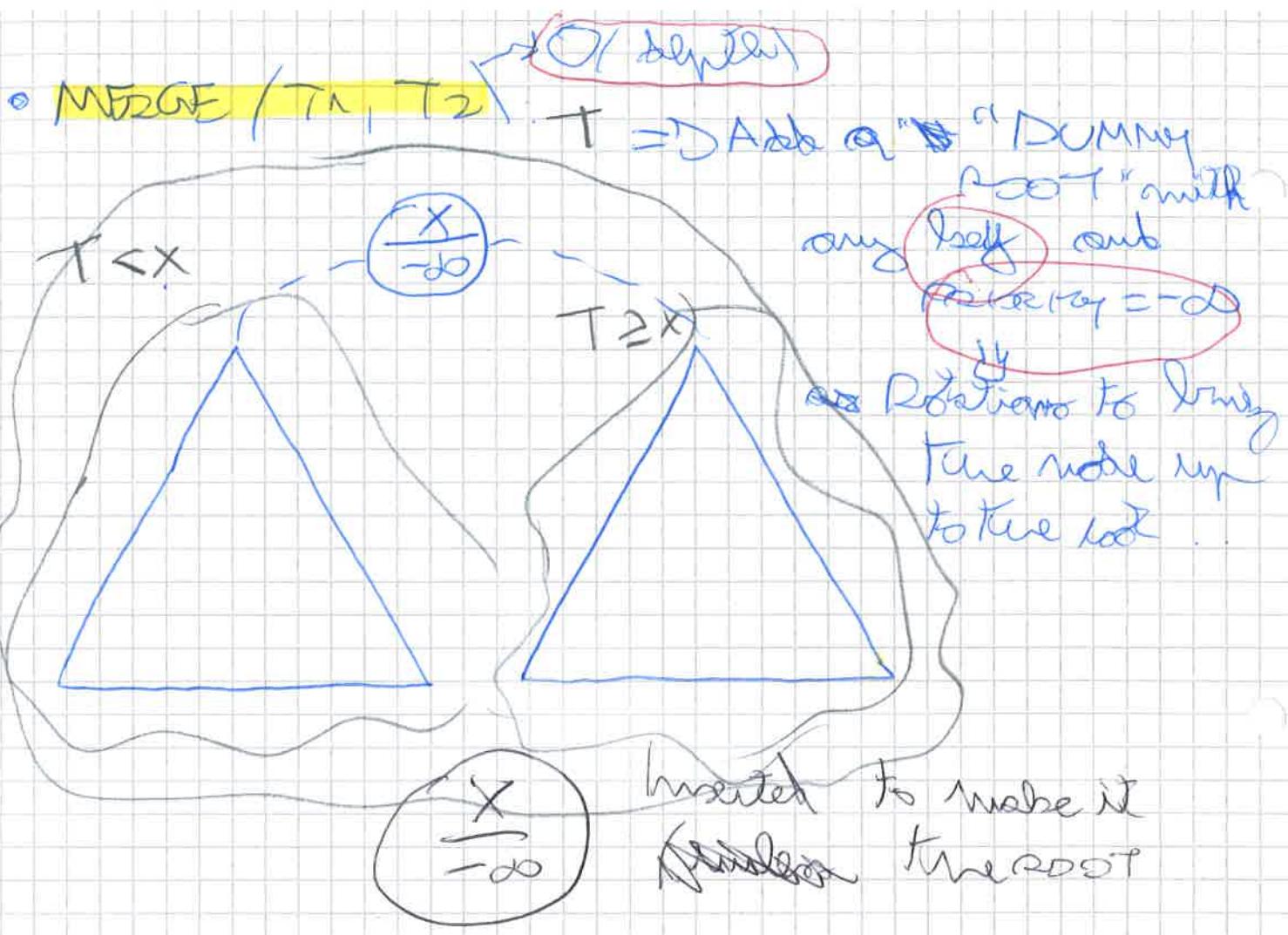
Root

\Rightarrow Push the node up to make it ~~a leaf~~

- If $K \notin \text{Tree}$:

\Rightarrow Insert (\underline{K}) , then

perform SPLIT when it's at the root.



RANDOMIZED TREAP.

A RANDOMIZED TREAP is a TREAP, where the priorities are assigned randomly to nodes uniformly and independently.

Since priorities are random, each node is equally likely to have the SMALLEST priority.

FACT:

If the PRIORITYES are picked randomly, then the TREAP is BALANCED $\Rightarrow \text{DEPTH} = O(\log n)$

In OPS. & insertion, deletion, search take $O(\log n)$ time.
DEPTH!

PROOF: (like RAV. ODESSER)

KEYS: X_1, X_2, \dots, X_K = Nodes with the K^{th} -smallest increasing order key.

We now define an INDICATOR RV:

$A_{(K)}^i$ = 1 iff. X_i is an ANCESTOR of X_K .

$$\text{depth}(X_K) = \sum_{i=1}^n A_{(K)}^i$$

(#ANCESTORS preceding node)
 K

$$\text{Take } E\{ \cdot \} = E\{ \circ \}$$

$$E\{ \text{depth}(X_K) \} = \sum_{i=1}^n E\{ A_{(K)}^i \}$$

$$A_{(K)}^i = \begin{cases} 1 & \text{iff } X_i \text{ is ANCESTOR of } X_K \\ 0 & \text{Otherwise} \end{cases}$$

$$\begin{aligned}
 E\{ \text{depth}(X_K) \} &= \sum_{i=1}^n E\{ A_K^{(i)} \} = \\
 &= \sum_{i=1}^n 1 \cdot P\{ A_K^{(i)} = 1 \} + 0 \cdot P\{ A_K^{(i)} = 0 \} \\
 &= \sum_{i=1}^n P\{ A_K^{(i)} = 1 \} \quad (\text{Given two keys, when } \\
 &\quad \text{one key is an ancestor of the other one } X_K)
 \end{aligned}$$

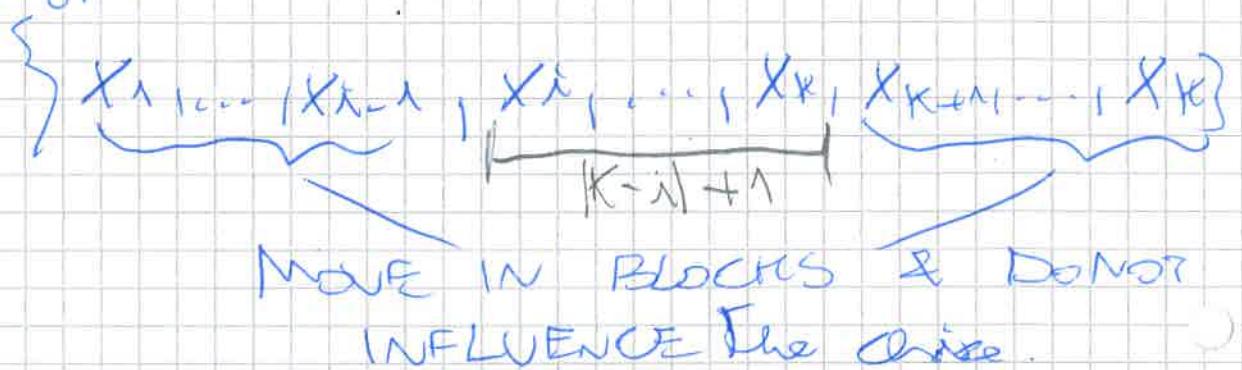
OBSERVATION:

Let's take: $\{X_1, \dots, X_{i-1}, X_i, \dots, X_K, X_{K+1}, \dots, X_n\}$

~~(X_{i+1}, ..., X_{K-1}, X_K, ..., X_n)~~

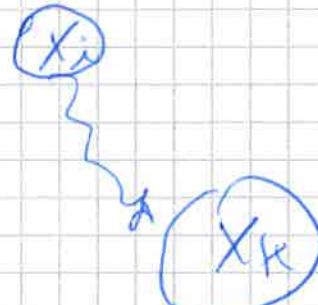
Given that priorities are assigned randomly, when no X_i is an ANCESTOR of X_K ?

Firstly, we notice that:



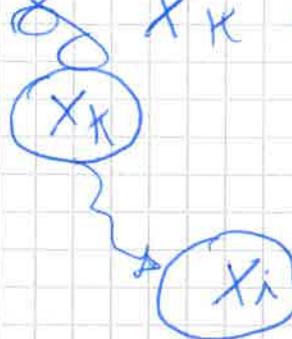
⇒ We have 3 cases (we will pick the one where X_i is an ancestor of X_K)

④ Priority of X_i is the MINIMUM one.



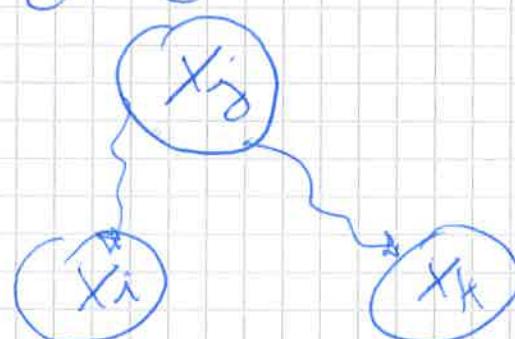
⇒ X_i is the root & surely an ancestor of X_K .

② Priority of x_k is the MINIMUM one.



x_k is the root & an ancestor of x_i \Rightarrow

③ Priority of x_i is the MINIMUM one.



$\Rightarrow x_i$ can't be an ancestor of x_k , as x_i and x_k have gone to two different branches.

\Rightarrow We are interested in SITUATION ①.

What is the P. of last case false verified?

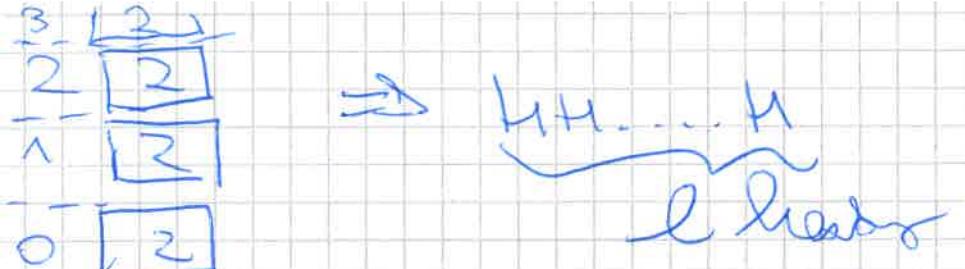
$$\sum_{i=1}^n P\{A_k^i = N\} = \sum_{i=1}^n \frac{1}{\sqrt{k-i+1}} = O(\log n)$$

Q.E.D.

$P\{PRIORITY(x_i) \neq PRIORITY(x_k)\}$

$x_1 x_2 \dots x_i x_{i+1} \dots x_{k-1} x_k x_{k+1} \dots x_n$

\Rightarrow We have hence proved that TREES with random priority are balanced & their depth is $O(\log n)$



\Rightarrow take at level (ℓ) , we need to have
at least l HEADS!

$$= \sum_{i=1}^n i \cdot (\text{Pr}[H_i])^\ell = \sum_{i=1}^n i \cdot \left(\frac{1}{2}\right)^\ell$$

$$= \left(\frac{1}{2}\right)^\ell \cdot \sum_{i=1}^n i = \cancel{\sum_{i=1}^n i} \cdot \frac{n}{2}^\ell \quad (\text{If } \ell = \log n)$$

\rightarrow height

$$\lg \ell = c \cdot \log n$$

$$\text{Pr}[L \geq c \cdot \log n] \leq \frac{n}{2}^\ell = \frac{n}{2^{c \cdot \log n}}$$

WEIGHT OF

$$\text{SKIP LIST} = \frac{n}{\cancel{\text{height}}} = \frac{n}{nc}$$

$$= \frac{1}{n^{c-1}} = \frac{1}{n}$$

$$\lg c = 2$$

$$\Rightarrow \text{Pr}[L \geq 2 \cdot \log n] \leq \frac{1}{n}$$

with
high
Pr[B!]

\Rightarrow Decreasing by $\frac{1}{n}$ ratio!

\Rightarrow We have proved that with high P_j ,
a skip list has $O(\log n)$ LEVELS

$E\{\#LEVELS\} = O(\log n)$ mitte hoog prob.

SEARCH ($K \rightarrow$) Start from the TOP-MOST LIST at the LEFT-MOST CORNER.

→ Go as far as possible to the right, then go down [and vice versa].
[Going down, you increase the "depth"]
Ex: $\frac{n}{k}$ $\leq K$

$$+ \text{next}(n) \leq K$$

\Rightarrow Go right

$$\text{if next}(n) \rightarrow R$$

\Rightarrow Go down

, based on the #LEVELS

TIME:

$O(\log n)$

#LEVELS $\approx O(\log n)$ mitte hoge prob.

PROOF:

We need to prove that the depth (max. #LEVELS) is logarithmic to the #keys.

~~P{#LEVELS in skip list $\geq l$ } $\leq \frac{1}{2^l}$ (P3) (A2)~~

MAX. #LEVELS

$$P\left\{\begin{array}{l} \text{MAX. } L \geq l \\ \text{in skip list} \end{array}\right\} \leq \sum_{i=1}^n P\{L/i \leq l\}$$

level of
item i

Ex: To reach level 3 at the item i , we have first at least 3 levels (and possibly 3 card).



$\Rightarrow \underbrace{111\dots1}_{l \text{ levels}}$

\Rightarrow take at level (l) , we need to have
at least l HEADS!

$$= \sum_{i=1}^n \Pr\{H_i\}^l = \sum_{i=1}^n 1 \cdot \left(\frac{1}{2}\right)^l$$

$$= \left(\frac{1}{2}\right)^l \cdot \sum_{i=1}^n 1 = \cancel{\sum_{i=1}^n} \frac{n}{2} l \quad (\text{if } l = \log n) \\ \text{+ height}$$

$$\lg l = c \cdot \log n \quad P\{L \geq c \cdot \log n\} \leq \frac{n}{2^l} = \frac{n}{2^{c \cdot \log n}}$$

$$\text{WEIGHT of} \\ \text{skip list} = \frac{n}{\cancel{2^{c \cdot \log n}}} = \frac{n}{n^c}$$

$$= \frac{1}{n^{c-1}} \approx \sum \frac{1}{n}$$

$$\lg c = 2$$

$$\Rightarrow P\{L \geq 2 \cdot \log n\} \leq \frac{1}{n} \quad \text{mitter fach Prob!}$$

\Rightarrow Decreasing by $\frac{1}{n}$ ratio!

\Rightarrow We have proved that mitter fach P_j
R-Skip List has $O(\log n)$ LEVELS.

COMPLEXITY OF SEARCHING

PROOF:

(and deleting / inserting)
an element

$$\text{SEARCH} = \# \text{VERTICAL STEPS} + \\ 2 \text{ UP} + \# \text{LEFT STEPS} \quad \# \text{HORIZONTAL STEPS}$$

$$E\{\text{Search}\} = E\{\# \text{VERTICAL STEPS}\} + \\ E\{\# \text{HORIZONTAL STEPS}\}$$

We know that:

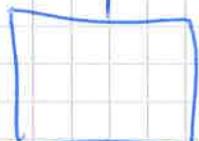
$$E\{\# \text{VERTICAL STEPS}\} = L = O(\log n)$$

(Height of the skip list)

~~We know that.~~

~~P} Go up from a node~~
~~This is a vertical link~~

$$P\{U\} = \frac{1}{2}$$



$$P\{ \text{Go up from a node} \} = P\{ \exists \text{ VERTICAL LINK} \} \text{ from a node}$$

$$= P\{U\} = \frac{1}{2} \Rightarrow E\{\# \text{VERTICAL LINKS}\} "O(\log n)"$$

P} Go left from a node

= P{ ~~Vertical Link from a node~~ }

$$= P\{T\} = \frac{1}{2} \Rightarrow E\{\# \text{HORIZONTAL LINKS}\} "O(\log n)"$$

~~E\{\# Nodes that have NO VERTICAL LINK~~

~~# Nodes that do have a VERTICAL LINK~~

$$E\{S_{\text{SEARCH}}\} = E\{\# \text{VERTICAL STEPS}\} + \\ E\{\# \text{HORIZONTAL STEPS}\} \\ - \text{O}(1 \text{ loop}) \quad \text{Q.E.D.}$$

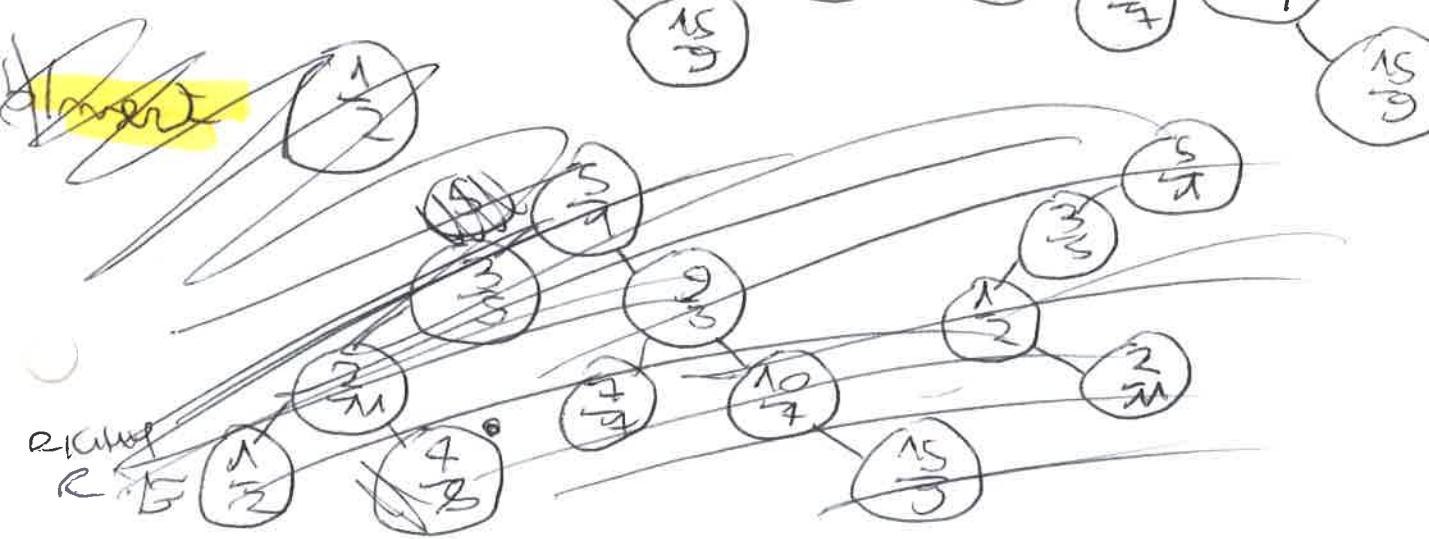
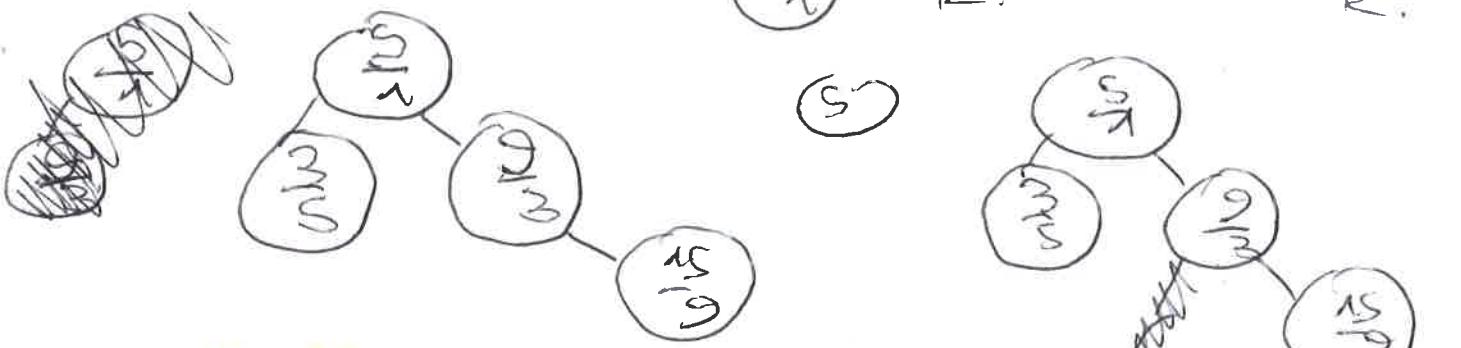
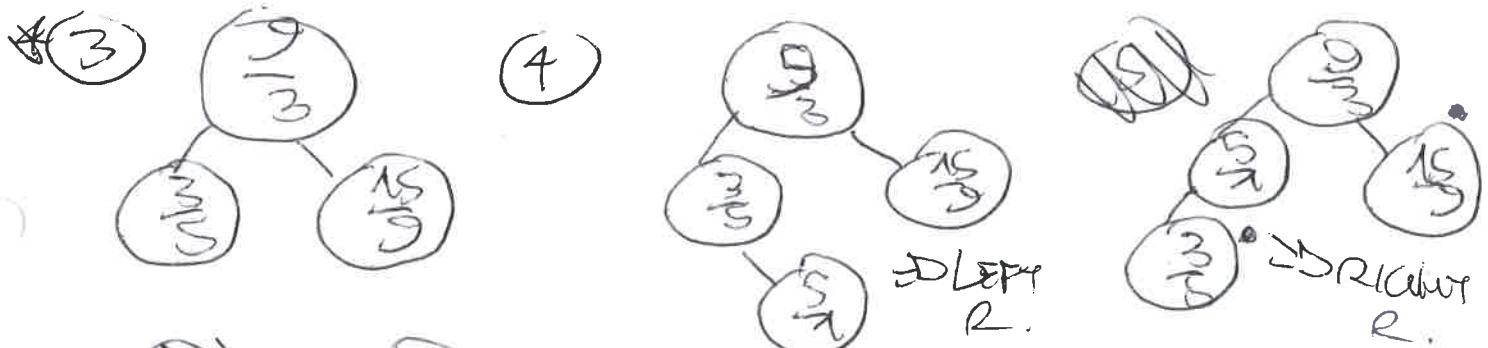
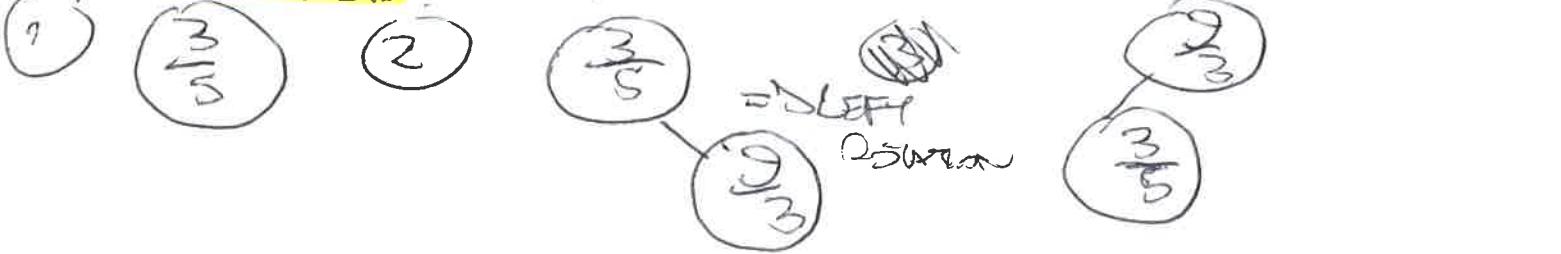
As the process of going up

↑
The process of going left

TREAP - CLASS EXERCISE:

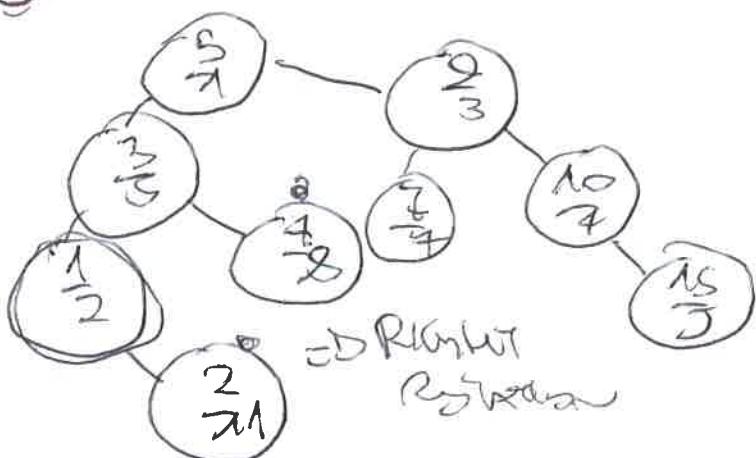
~~S = { (3, 5), (9, 1), (7, 15), (5, 1), (10, 2), (14, 1), (2, 14), (1, 15) }~~

① BUILD TREAP



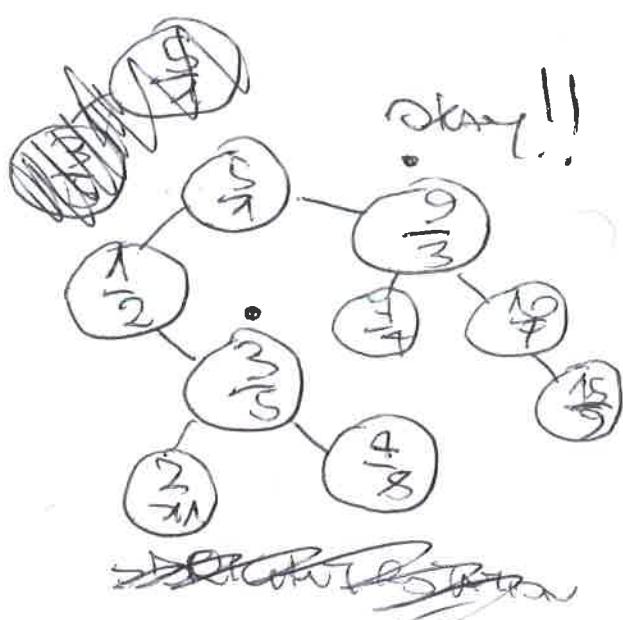


\Rightarrow RIGHT Rotation

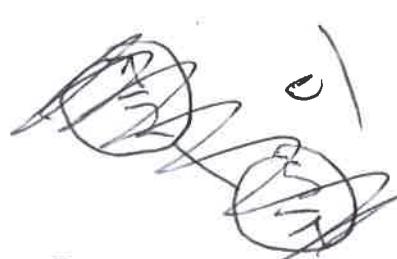


\Rightarrow RIGHT Rotation

IDEA: Perform your LEFT RIGHT ROTATION, then invert the remaining items (leaves) in the right position of the tree.



\Rightarrow Skew Tree



c) Needs to perform a SPLIT at $\underline{6}$.
is inserted, and then "brought up".



RIGHT

\Rightarrow Rotation

