

Algorithm: Communication of the ACM

Def. of algorithm (by Knuth)

Finite no. of steps

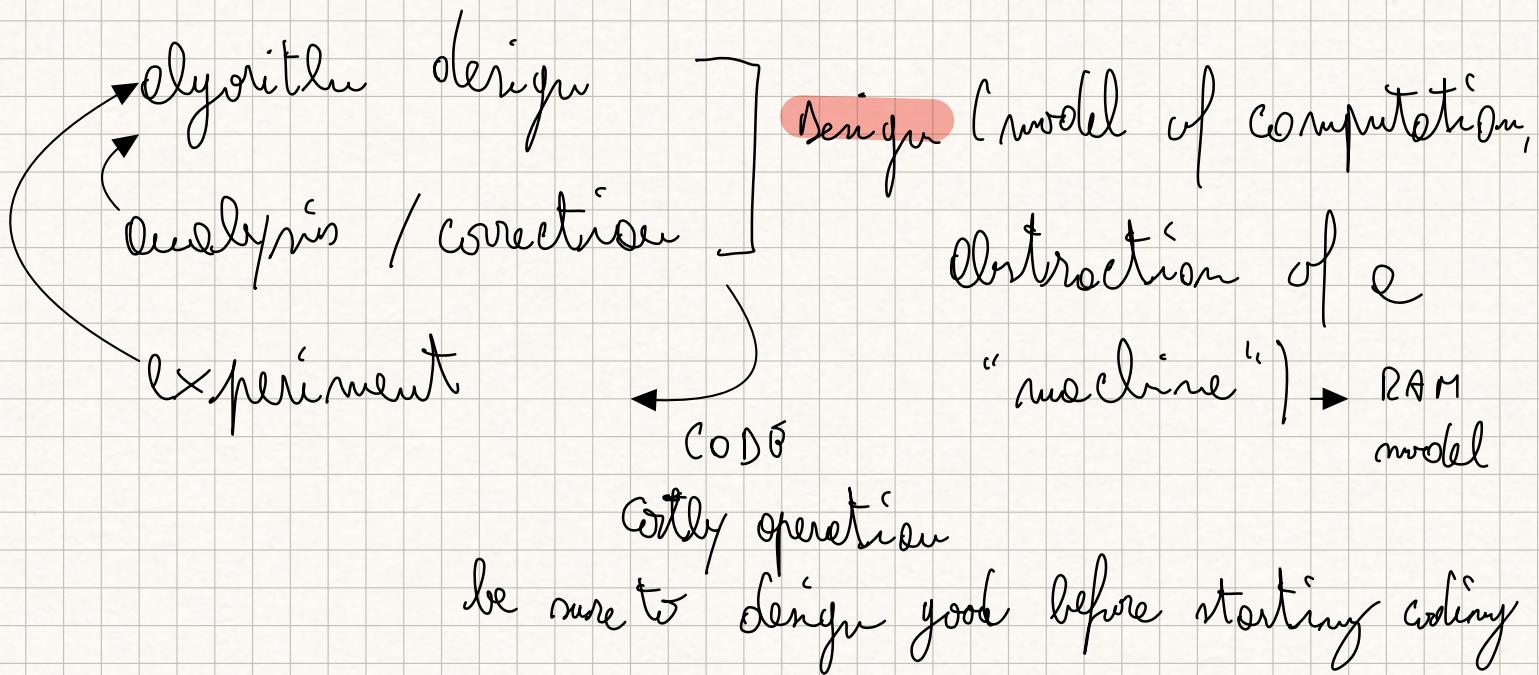
W Prop: • DEFINITION = unambiguous \rightarrow programming lang.

• EFFECTIVE = $\{ \text{input} \rightarrow \text{output} \}$

• FINITE $\rightarrow O(1)$ TIME

• INPUT and OUTPUT

problem abstraction / formulation

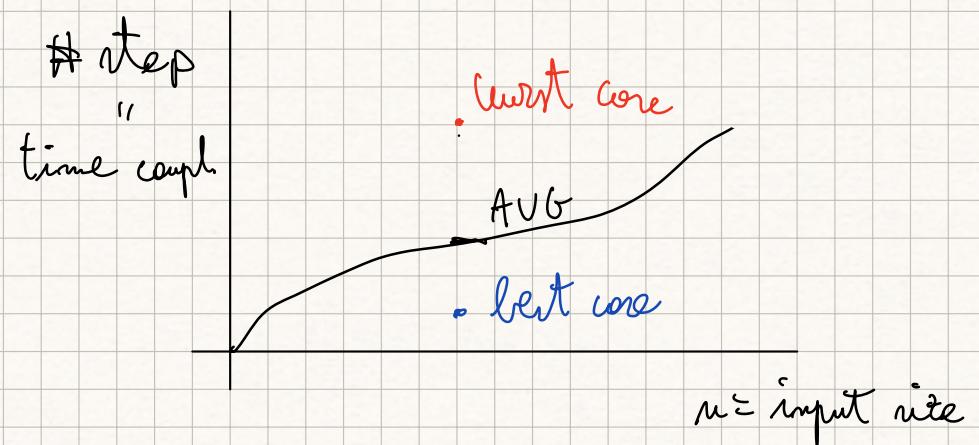


• time complexity of an alg. $A \approx \# \text{steps}$ (comparison)

$g(\text{Input})$

↳ considers only input

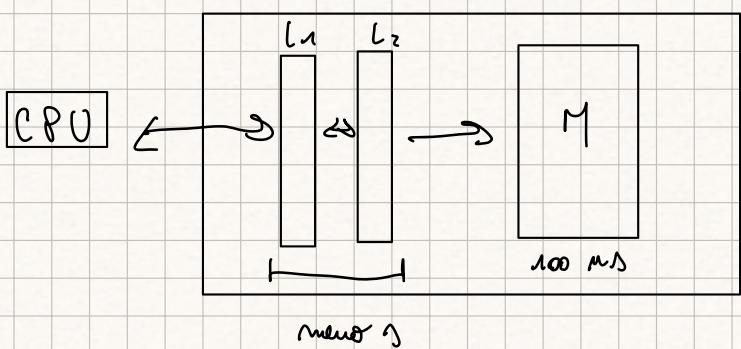
the min. does it worst case



• to compare alg. A_1 and $A_2 \rightarrow \# \text{steps}$

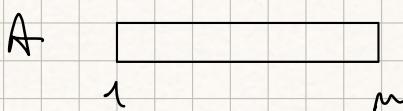
for BIG DATA $n \rightarrow \infty$

• memory \rightarrow hierarchy of memory levels



→ extremely complex. We need approximation.
Ex. external memory model

Do :



$\text{sum}(A) : \text{Elap} > 0$

for $i=1$ to m

$\text{temp} \leftarrow A[i] / O(n)$

PAM. oh ely.

A S, b

Jump ↳ block size within the Page B

$A_j = A[j \cdot b + 1, (j+1) \cdot b]$ → logical block

$A_0 = A[1, b]$

$A_1 = A[b+1, 2b]$

if $S=1 \rightarrow$ reading

• ref of block I occurs

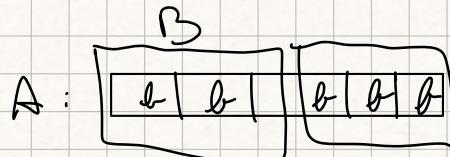
$$j = S \cdot i \bmod \frac{m}{b}$$

$i = 0, 1, \dots, m/b - 1$

$S > 1 \rightarrow$ different pattern of memory access

↳ co-prime with $\frac{m}{b}$ (per every n rows the bw occurs at t th block)
 \downarrow
M.C.D. = 1 item)

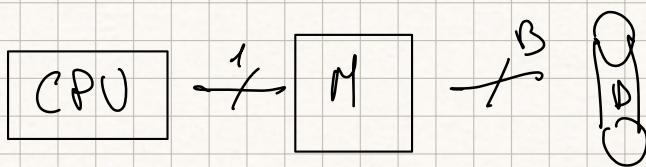
number of max pass $\rightarrow \frac{B}{b}$



$$O(I/O) = \min \left\{ S, \frac{B}{b} \right\} \cdot \frac{m}{B}$$

$$O(\text{time}) = n$$

Two-level memory model



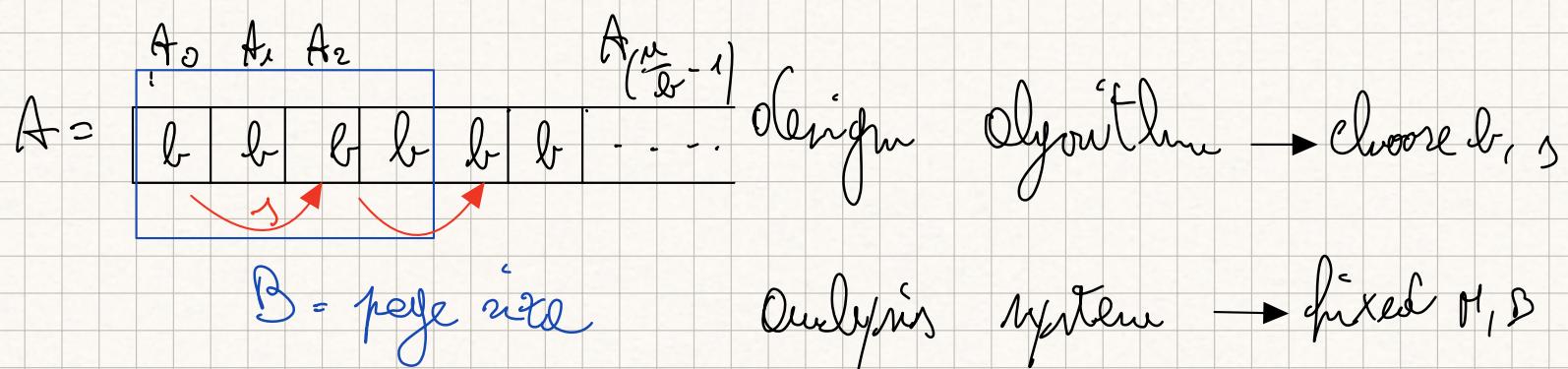
internal
memory /
cache

disk / internal memory

HDD
SSD

}

Two level memory
depends on how much
data is here



$$A_j \text{ where } j = i \cdot s \text{ and } \frac{m}{B}$$

$$0, 1, 2, \dots, m/B - 1$$

Time complexity $= O(n)$

$$I/O \text{ complexity} = O \left(\min \left\{ \frac{B}{b}, s \right\} \cdot \frac{m}{B} \right) =$$

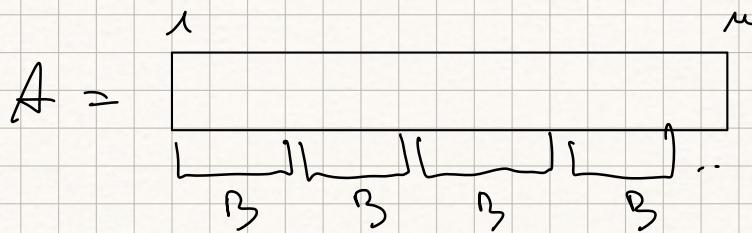
$$s=1 \rightarrow \frac{m}{B} \text{ Scanning}$$

$$s>1 \rightarrow s \cdot \frac{m}{B}$$

$$s < \frac{B}{b}$$

Searching in a sorted array

Key $\in A$?



binary search : comparison based algorithm

$$\text{Time} = O(\log_2 n)$$

$$I/O = \log_2 n - \log_2 B = \log_2 \frac{n}{B}$$

v v
 ceiling floor
 regime regime

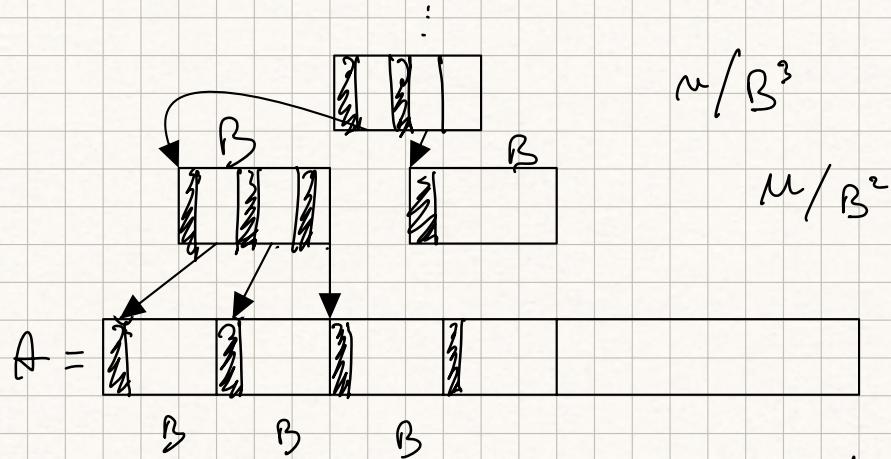
$B = 32 \text{ Kbytes}$

- ↖ seek time $5\mu s$
- ↙ fetch time

8 byte int's

$$\# \text{ items} \times \text{page} = \frac{2^{15}}{2^3} = 2^{12} \approx 4000 \text{ items} \log_2 \frac{2^{30}}{2^{12}} \approx \frac{30}{12} = 18 I/O$$

B-Tree :



n/B

the root will be $\frac{n}{B^e} \leq 1 \rightarrow \log_B n \leq l$

$$\text{I/O complexity} = \log_B n = \frac{\log_2 n}{\log_2 B}$$

$$\frac{\log_2^{2^{30}}}{\log_2^{2^{11}}} = \frac{30}{11} \approx 3 \text{ I/O much better than before}$$

$$\log_2 \frac{n}{B} \Leftrightarrow \log_B n$$

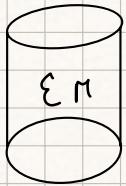
(3)

B-tree is static, meaning that it can move inside it with arithmetic (no pointer)

Spatial locality: i use items inside the same block

Temporal locality: i use items during the execution items of the same blocks

M



step

computation $1 - \epsilon = 70\%$

cost = 1

internal memory

external memory

$$M = (1 + \epsilon) M$$

$\epsilon = 30\%$

memory access

internal

1

cost = 1

c

external

prob. of finding needed

an item in the external memory

$$P(\epsilon) : \frac{\epsilon \cdot M}{M} = \frac{\epsilon}{1 + \epsilon}$$

Copy-Cost

$$\text{Copy cost step} : (1 - \epsilon) \cdot 1 + \epsilon [(1 - P(\epsilon)) \cdot 1 + P(\epsilon) \cdot C]$$

internal mem. access

→ all the data

↓
external memory access

$$< 0 + 0 \cdot P(\epsilon) \cdot C \approx 0.3 + 0.3 \cdot \frac{\epsilon}{1 + \epsilon} \cdot 10^6 \approx 2 \cdot 10^5 \cdot \frac{\epsilon}{1 + \epsilon} + \text{cost}$$

$\frac{\epsilon}{1 + \epsilon}$ $P(\epsilon)$

$$P(\epsilon) = \frac{1}{1000}$$

occurring to the

external memory could cost

o Let blue with low probability $P(\epsilon)$

Sorting

RAM

permuting

M

working

n log n

SORTING

$$A_{\pi} = \boxed{C \ A \ B \ D} \quad \text{for } i=1 \text{ to } n$$

$$A = \boxed{A[B \ C \ D]} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{PERMUTATION}$$

$$A_{\pi}[i] = A[\pi[i]]$$

- Compute the sorted permutation π
 - Implement $\pi[m]$
- $m \cdot \log m$

| $O(\cdot)$ | RAM | 2-level (external memory) |
|--------------|------------|--|
| permuting | m | $\min \left\{ \frac{m}{B} + \text{sort}, \frac{m}{2} + \text{sort} \right\}$ |
| All | $n \log m$ | $I/O_s \leq 3 \text{ read } O\left(\frac{m}{B}\right)$ |

Rephrase the permuting problem as a sorting problem

Memory \rightarrow File of tuples
sequence

array

$$1) \quad \left\{ \begin{array}{l} S = (A, 1), (B, 2), (C, 3), (D, 4) \quad \text{cost: } O\left(\frac{m}{B}\right) I/O, \\ \text{item} \quad \text{position} \end{array} \right.$$

$\pi = \langle 3, 1 \rangle \quad \langle 1, 2 \rangle \subset \langle 2, 3 \rangle \subset \langle 4, 4 \rangle$
 source destination
 cost: $O\left(\frac{n}{B}\right)$ I/Os

2) Sort π by first element: $\langle 1, 2 \rangle \langle 2, 3 \rangle \langle 3, 1 \rangle \langle 4, 4 \rangle$

3) Scan π and ς : $\langle A, 2 \rangle \subset \langle B, 3 \rangle \subset \langle C, 1 \rangle \subset \langle D, 4 \rangle$

4) Sort π by second component $\langle C, 1 \rangle \subset \langle A, 2 \rangle \subset \langle B, 3 \rangle \subset \langle D, 4 \rangle$

$\text{Sort}(n, n, B) = O\left(\frac{n}{B} \cdot \log \frac{n}{B}\right)$ I/Os
 elements memory page size

$$n > 16 \text{ Gb} \approx 2^{34}$$

$$\frac{n}{B} = 2^{19}$$

$$B = 32 \text{ Kb} = 2^{15}$$

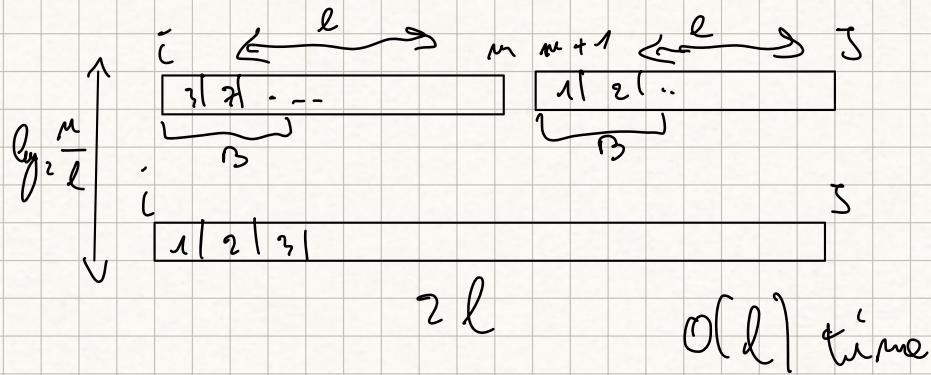
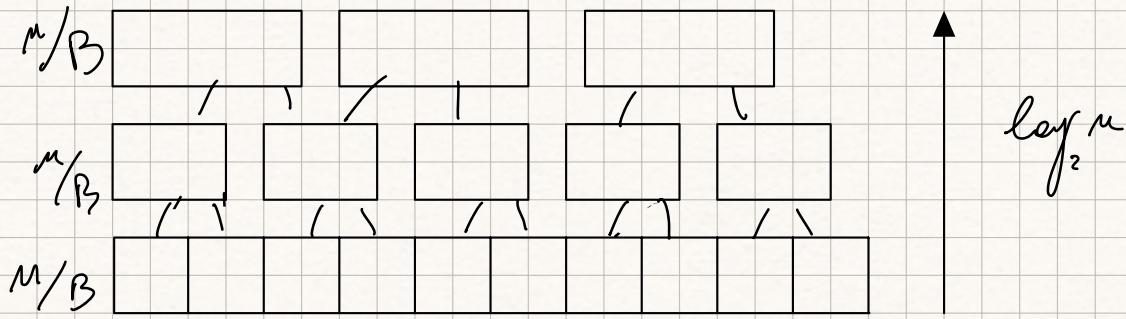
$$\log \frac{n}{B} = \frac{\log_2 \frac{n}{B}}{\log_2 \frac{B}{B}} = L$$

$$n = 2^{60} \quad \# \text{Scans} = \frac{\log_2 \frac{2^{60}}{2^{34}}}{\log_2 2^{19}} = \frac{26}{19} \approx 2$$

When permute is better than sort!

$$\frac{n}{B} \cdot L \leq n \rightarrow B \geq L \quad L \text{ is } 2 \text{ w permutations and writing are equal}$$

Binary Mergesort:

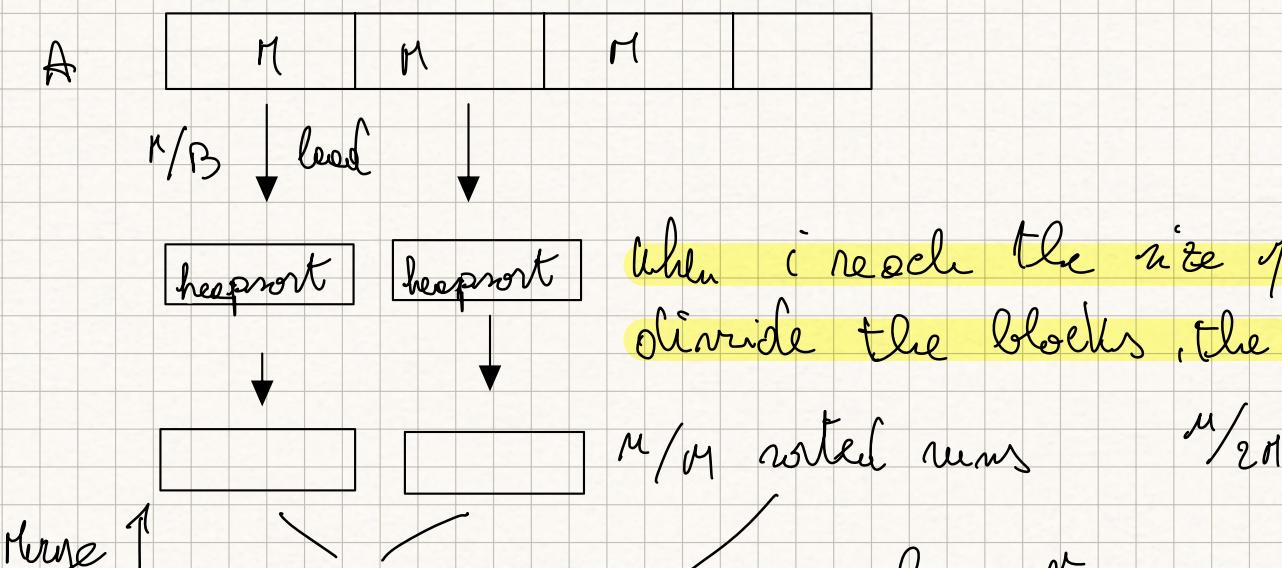


$$\text{Binary MS} = O\left(\frac{m}{B} \log_2 n\right) \text{ I/O}$$

$$\text{Considering Memory} = O\left(\frac{m}{B} + \frac{m}{B} (\log_2 n - \log_2 B)\right) = O\left(\cancel{\frac{m}{B}} + \frac{m}{B} \left(\log_2 \frac{m}{B}\right)\right)$$

$\ell \leq B$ $\ell > B$

dim. of the partition ℓ len or eq. to block size

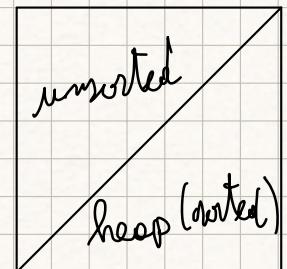


$$\log_2 \frac{m}{2^m}$$

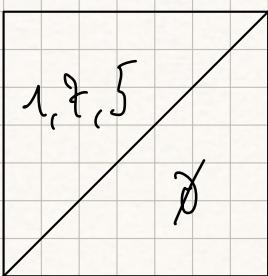
Snow Plow:

$$A = [1, 2, 5] \mid [3, 2, 6] \mid [4, 3, \dots]$$

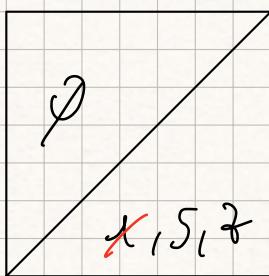
$$m = 3$$



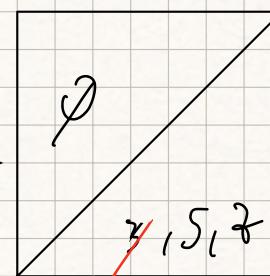
phase:



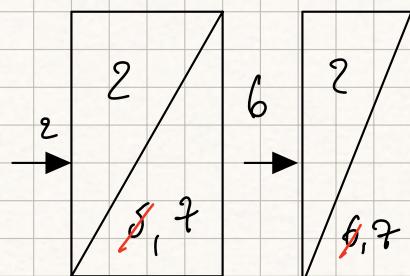
work



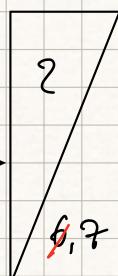
next
3



2



6

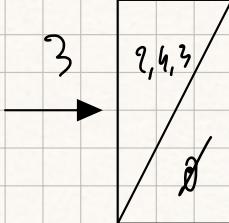
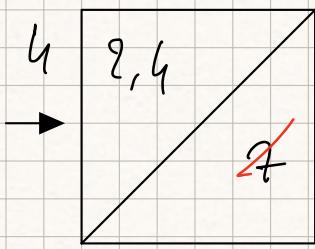


1 if $3 < 1 \rightarrow 0$ 3

5 6

if $3 \geq 1 \rightarrow 4$

1 3 5 6 7

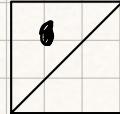


7

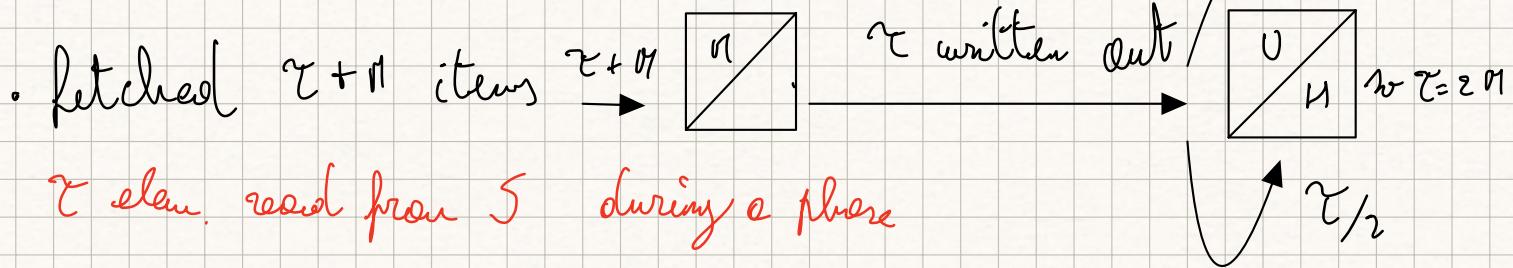
The expected length of a run created by SnowPlow is

$2m$.

- load m items in memory



- load τ items until the heap is empty and (unsorted) \cup
if full again $P(\text{item goes to } M) = 1/2$



Qs: $M = 2$ $1, 8, 3, 2, 5, 0, 4, 6$ 3 runs

$1, 3, 8$

$2, 5$

$0, 4, 6$

1) **Phase**: load M unsorted elements in the min-heap H

2) **Run**: writes to the output the minimum item within H ,
and load the next item from S in an auxiliary array U
located in the memory, if $\text{next} < \text{current}$; otherwise load
next item in H if $\text{next} > \text{current min}$

3) **Last of a Phase**: H is empty and U consists of M unsorted
elements. At the end you only merge the runs, this time
from the simplified method of mergesort.

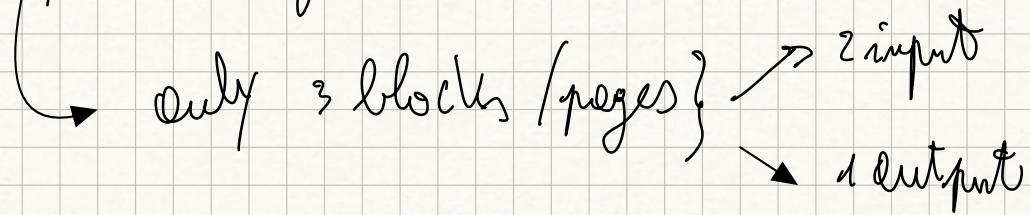
If I increase M to $2M$ I save 1 level.

$$\# \text{ levels} = \log_2 \frac{M}{M} \rightarrow \log_2 \frac{M}{2M} = \log_2 \frac{M}{M} - \log_2^2 = \log_2 \frac{M}{M} - 1$$

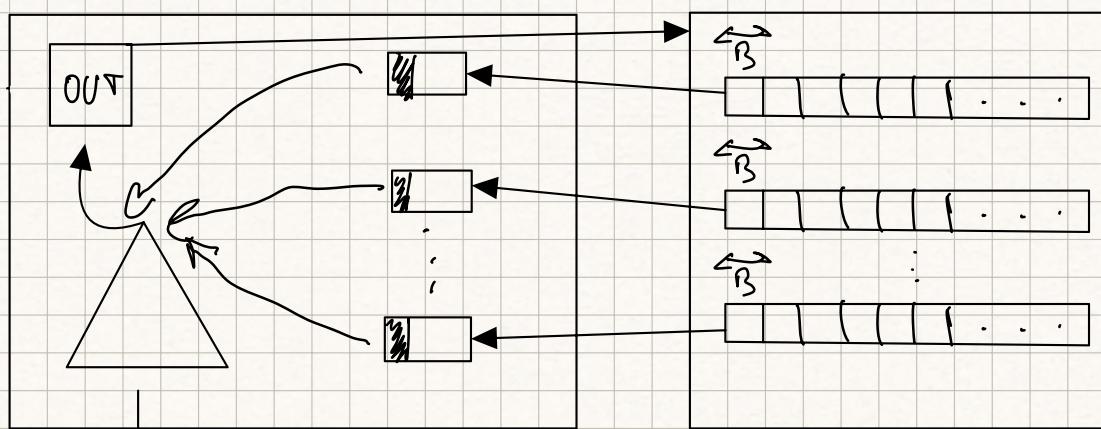
$$\text{Binary Merge} = O\left(\frac{m}{B} \log_2 \frac{m}{M}\right) \%$$

$\frac{m}{B}$

not fully exploiting the internal memory



Multi-way merge sort:



Run 1
Run 2
Run 3
Run 4

grouped

$\min - heap$
of size K → contains pair composed by $(\min elem, run)$.

$$O(K \log_2 K)$$

Rs:

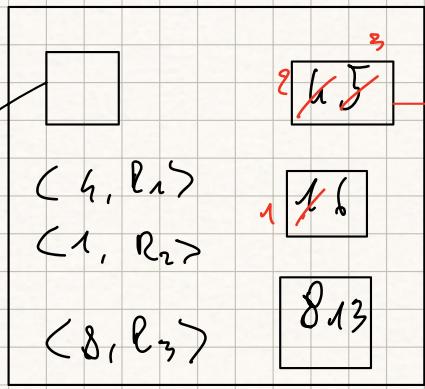
| | | | | |
|---------------|-------|-------|-----|-------|
| 15 | 10 11 | 15 21 | ... | R_1 |
|---------------|-------|-------|-----|-------|

$$B = 2$$

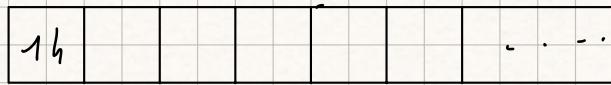
| | | | | |
|----|-----|-------|-----|-------|
| 16 | 7 9 | 14 18 | ... | R_2 |
|----|-----|-------|-----|-------|

| | | | | |
|------|-------|-------|-----|-------|
| 8 13 | 20 25 | 30 46 | ... | R_3 |
|------|-------|-------|-----|-------|

M~~O~~ M~~O~~ R Y



Load a new block 1011



How can i put all the elements in the memory?

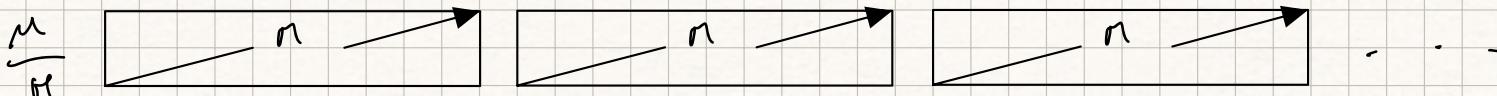
$$(k+1)B + \underbrace{\text{heapsize}}_{k} \leq M \quad k \text{ runs } + B \text{ for merging}$$

$$(k+1)B + \cancel{k} < M$$

$$(k+1) < \frac{M}{B}$$

$$k < \frac{m}{B} - 1 = O\left(\frac{m}{B}\right)$$

$\hookrightarrow 2^{15}$



$$\dots \sqrt{\left(\frac{M}{B} + O + \frac{M}{B}\right) \cdot \frac{M}{M}} = O\left(\frac{M}{B}\right)$$

$$O\left(\frac{M}{B}\right)$$

M/B

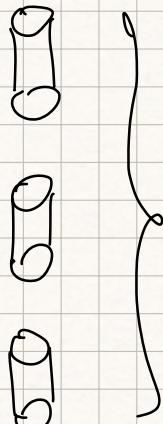
m/B

total cost: $O\left(\cancel{\frac{M}{B}} + \frac{m}{B} \log \frac{m}{B} \frac{M}{m}\right) I/O$,

initial number of runs

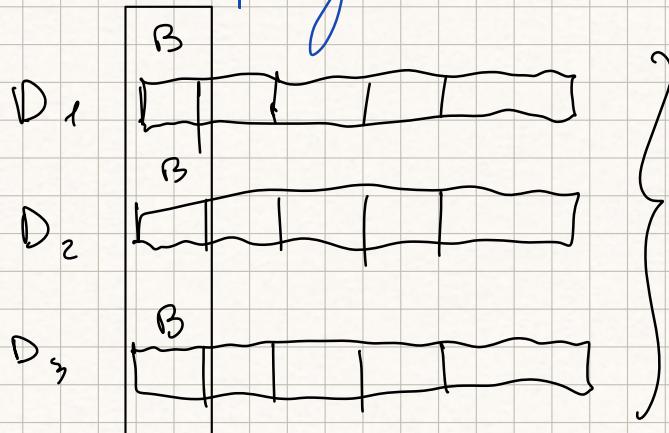
What happens when disks are more than 1? $D > 1$

parallelism; every disk can fetch one block

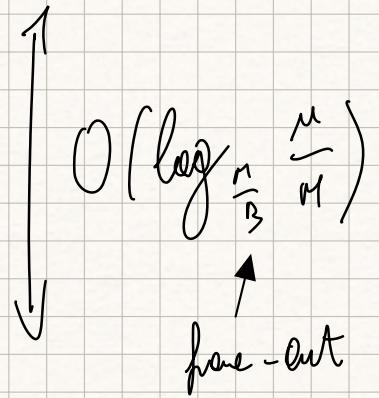


independent $\Rightarrow 1 I/O = 3B$ items
 $= 3$ blocks

Disk Striping: look at D disks as one single disk with $B' = DB$



$$D' = D_1 + D_2 + D_3$$



$$B' = D \cdot B$$

with $D=1$ $O\left(\frac{m}{B} \log_{\frac{m}{B}} \frac{m}{M}\right)$ $D=1$

↓ disk striping over D disks

$$O\left(\frac{m}{DB} \log_{\frac{m}{DB}} \frac{m}{M}\right) = O\left(\frac{m}{DB} \log_{\frac{m}{DB}} \frac{m}{M}\right) \text{ I/Os}$$

lower bound for writing: $\Omega\left(\frac{m}{DB} \log_{\frac{m}{DB}} \frac{m}{M}\right) \text{ I/Os}$

disk striping: $\frac{\frac{m}{DB} \log_{\frac{m}{DB}} \frac{m}{M}}{m}$

shaving factor =

optimized sh: $\frac{\frac{m}{DB} \log_{\frac{m}{DB}} \frac{m}{M}}{m}$

$$= \frac{\log_2 \frac{m}{M}}{\log_2 \frac{m}{DB}} \cdot \frac{\log_2 \frac{m}{B}}{\log_2 \frac{m}{M}} = \frac{\log_2 \frac{1}{B}}{\log_2 \frac{m}{B} - \log_2 D} = \frac{1}{1 - \log_{\frac{m}{B}} D} \geq 1$$

\downarrow
 $\frac{m}{B} \cdot \frac{1}{D}$

if $D \rightarrow \frac{m}{B}$ the shaving factor $\rightarrow +\infty$

disk striping is less efficient as the number of disks increases.

lower bound for permuting 5.2.2. per block.

lower bound for sorting ($D = 1$)

every node is a comparison

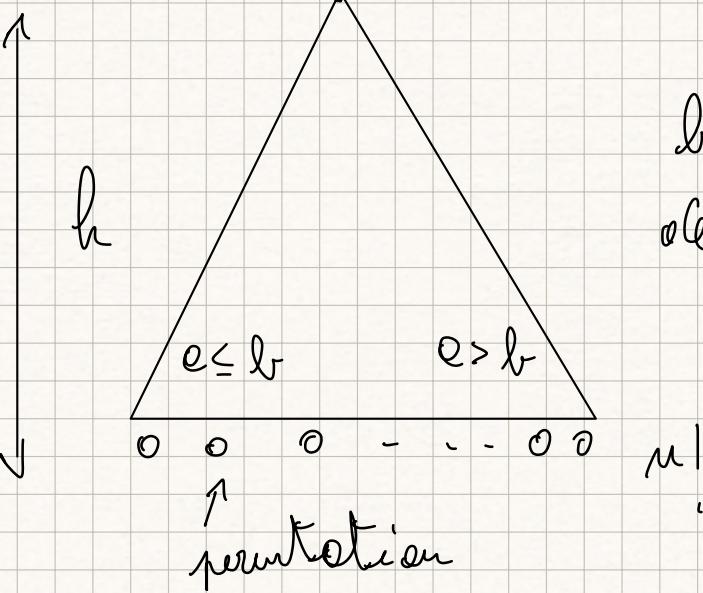
RAM only access

$$2^h \geq l = n!$$

$$h \geq \log_2 n!$$

$$\log_2 n!$$

$$n \log_2 n$$

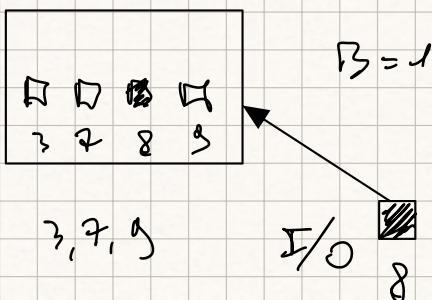


binary
selection
tree

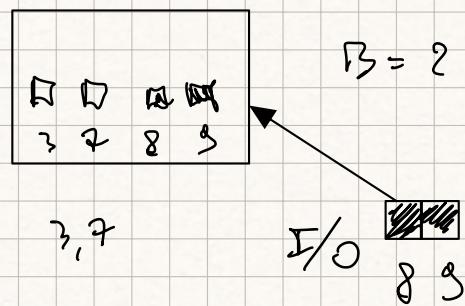
$$h \geq n \log_2 n \Rightarrow h = \Omega(n \log n)$$

RAM/DISK Access

$$n = u$$



$$n = u$$



max

$\binom{u}{2}$

is max

of possible allocation for $B=2$

① $\binom{u}{B}$ # ways B items yet distributed among $u-B$ items
(sorted) already in memory

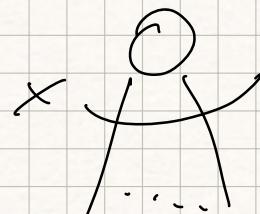
② first read \rightarrow count permutation of B items
new read
old read

next read \rightarrow no count

$t = \text{height of the tree}$

③

$$\frac{m}{B} \text{ new}$$



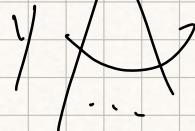
T/C

$$x = B! \cdot \left(\frac{m}{B}\right)$$

of steps that an algorithm has

to do

$$t - \frac{m}{B} \text{ old}$$



- old read

$$y = \left(\frac{m}{B}\right)$$

of nodes at the end

$$X^{\frac{m}{B}} \cdot y$$

$$\geq m!$$

$$m!$$

$$\left[B! \cdot \left(\frac{m}{B}\right) \right]^{\frac{m}{B}} \cdot \left[\left(\frac{m}{B}\right) \right]^{t - \frac{m}{B}}$$

$$(B!)^{\frac{m}{B}} \cdot \left(\frac{m}{B}\right)^t \geq m!$$

$$\text{property } \log_2 \alpha! = \alpha \log_2 \alpha$$

$$\frac{m}{B} \log_2 (B!) + t \log_2 \left(\frac{m}{B}\right) \geq \log_2 m!$$

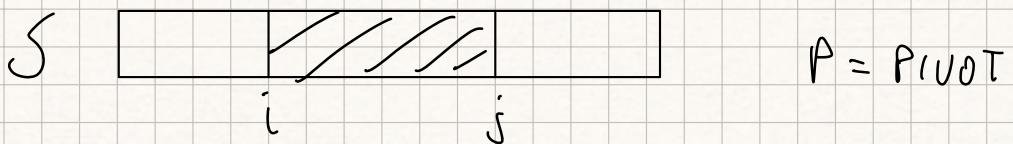
$$\log_2 \left(\frac{\alpha}{\beta}\right) = \alpha \log_2 \frac{\alpha}{\beta}$$

~~$$\frac{m}{B} \log_2 B + t \cdot B \log_2 \frac{m}{B} \geq m \log_2 m$$~~

$$t \cdot B \log_2 \frac{m}{B} \geq m \log_2 \frac{m}{B} \rightarrow t \geq \frac{m \log_2 \frac{m}{B}}{B \log_2 \frac{m}{B}} = \frac{m}{B} \cdot \log_2 \frac{m}{B}$$

Merge based \rightarrow Merge sort } Partition/divide
 Distribution based \rightarrow Quicksort } Divide and conquer \ Recursion / conquer
 Recombine

Quicksort : Partition (S, i, j)

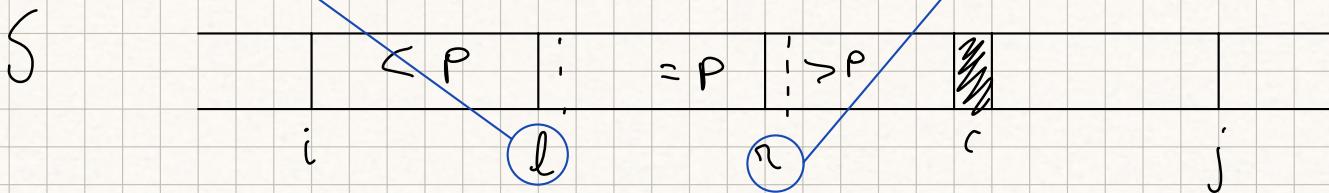


S $\boxed{<P} = P >P$ Three way partition

min elements = P

return
left return
right

min elements > P



$S[c] = P$ Swap $S[c]$ with $S[r]$, increment r e.c.

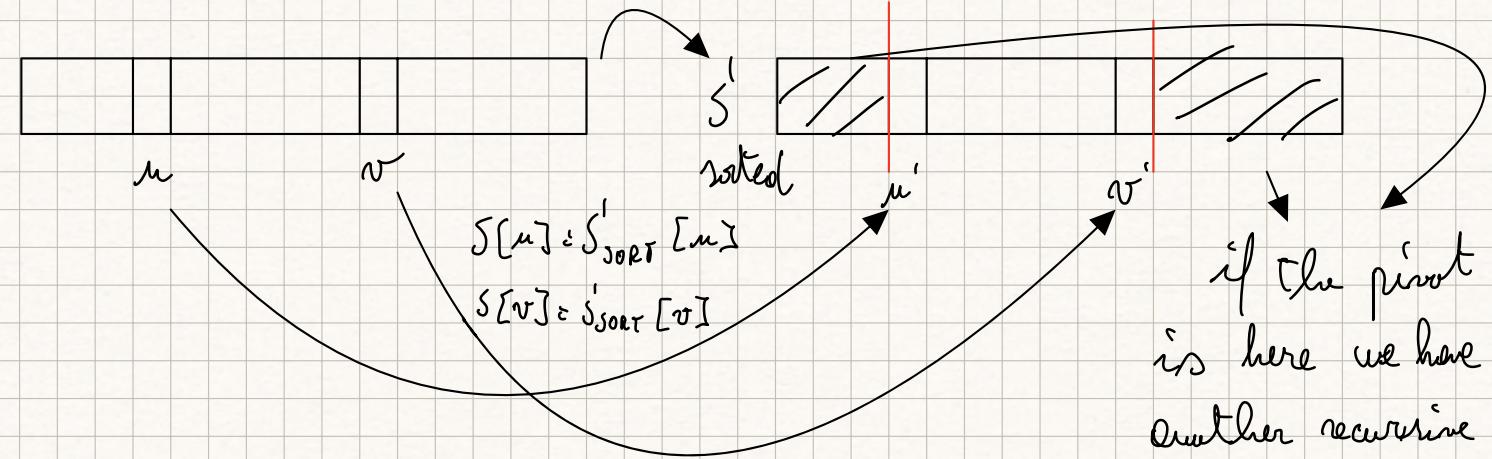
$S[c] < P$ Swap $S[c]$ with $S[l]$ & $S[c]$ with $S[r]$, inc. r, l, c

$S[c] > P$ increment c .

Expected time complexity of Quicksort : $O(n \log n)$

Random binary $X_{u,v}=1$ iff $S[u]$ compared with $S[v]$. One of them must be a pivot in order to be compared

$$\text{Exp \# comp.} = E \left[\sum_u \sum_{v>u} X_{u,v} \right] = \sum_u \sum_{v>u} E[X_{u,v}] = \sum_u \sum_{u>v} 1 \cdot P(X_{u,v}=1)$$



$P(m \text{ is pivot or } v \text{ is pivot})$

$$P(X_{m,v} = 1) = \frac{2}{n' - v' + 1} = \sum_{m'} \sum_{v' > m'} \frac{2}{n' - v' + 1} \leq 2 \log n$$

■

Bounded Quicksort (S, i, j):

while ($S[i] > m_0$) // parameter that stops recursion because i will work only in local memory
 $i = \text{select position of pivot}$

Swap $S[i]$ with $S[i]$

$p = \text{Partition}(S, i, j)$ // final position of the pivot

if $p \leq \frac{i+j}{2}$ { Bounded QS($S, i, p-1$); *faccio le ricorse nello resto perciò*
 $i = p+1$; } // remove tail recursion

else { Bounded QS($S, p+1, j$)

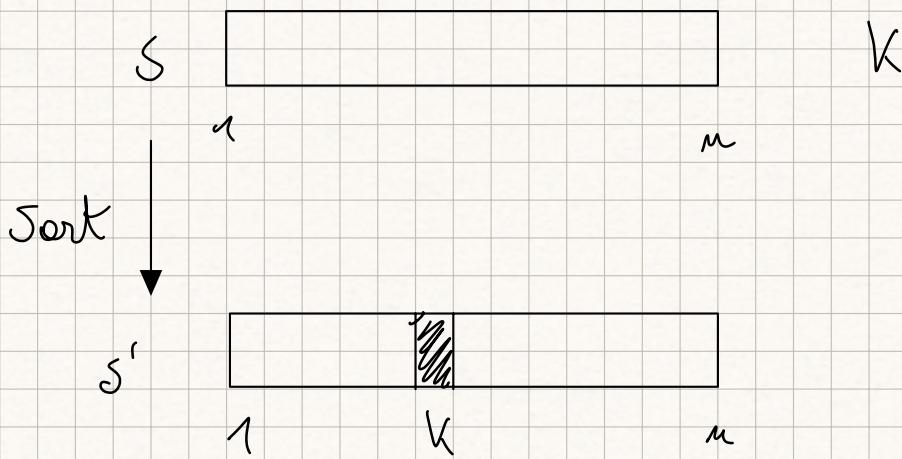
$j = p-1$;

}

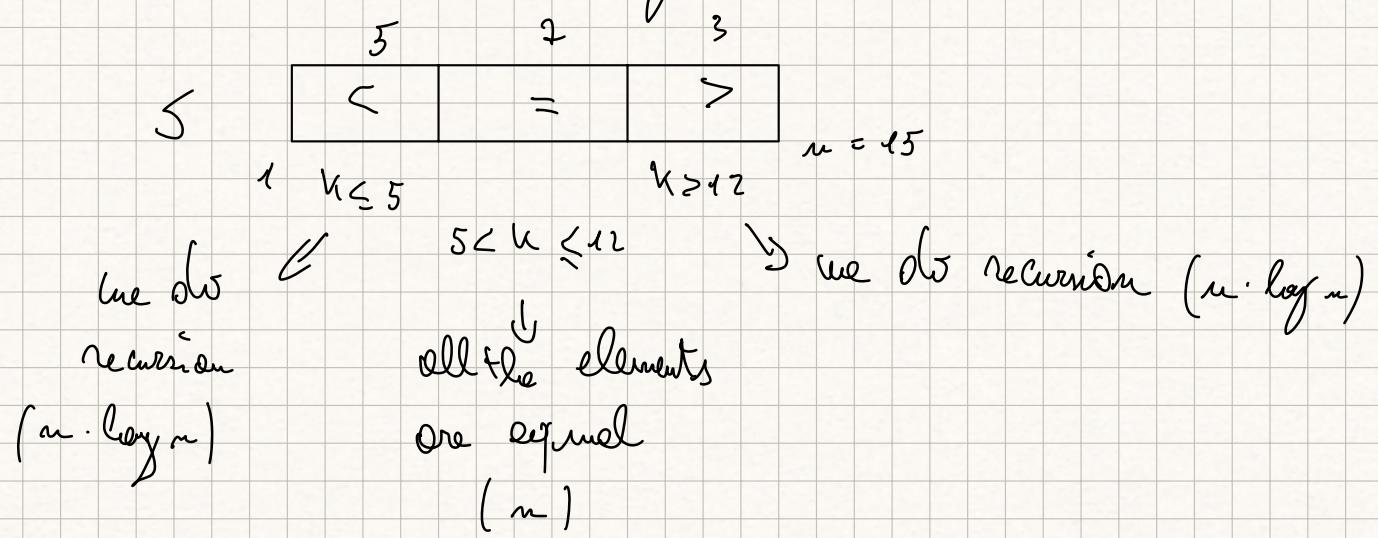
Insertion sort (S, i, j); // for very short arrays use IS.

→ better than heapsort

Select k -th ranked element from S (unordered)



We use 3-way partitioning



Random select (S, k): $\rightarrow k$ -th ranked item in S

$r = \text{random pos in } \{1, 2, \dots, m\}$

$S_r = \{x \in S, x < S[r]\}$

$S_> = \{x \in S, x = S[r]\}$

$$S_> = \{x \in S, x > S[s]\}$$

$$n_< = |S_<|$$

$$n_> = |S_>|$$

if ($k \leq n_<$) return Random Select ($S_<, k$)

else if ($k \leq n_< + n_>$) return S_{avg}

else return Random Select ($S_>, k - |S_<| - |S_>|$)

S_{sort}

| | | |
|--|--|--|
| | | |
|--|--|--|

$$\frac{n}{3} \quad \frac{n}{3} \quad \frac{n}{3}$$

$$T(n) = O(n) + p_c T(n_<) + p_s T(n_>) \quad \text{pivot should be here if chosen well}$$

Time Complexity

$$|S_<| < \frac{2}{3}n \quad |S_>| < \frac{2}{3}n$$

rest of the array

$$\frac{1}{3} T(n) \leq O(n) + \frac{1}{3} \underbrace{T\left(\frac{2}{3}n\right)}_{\text{good}} + \frac{2}{3} \underbrace{T\left(\frac{1}{3}n\right)}_{\text{bad}} \leq O(n) + \frac{1}{3} T\left(\frac{2}{3}n\right) + \frac{2}{3} \cdot \frac{1}{3} T(n)$$

$$\frac{1}{3} T(n) - \frac{2}{3} \frac{1}{3} T(n) \leq O(n) + \frac{1}{3} \frac{1}{3} T\left(\frac{2}{3}n\right)$$

worst case - when
k-th element is 1°
or last

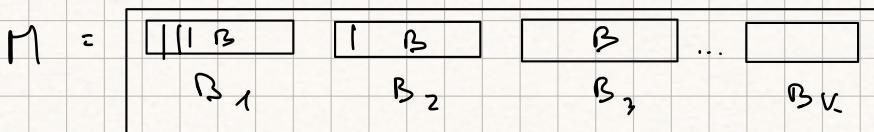
$$\frac{1}{3} T(n) \leq O(n) + \frac{1}{3} T\left(\frac{2}{3}n\right)$$

per il master theorem $\frac{1}{3} T(n) = O(n)$

In terms of I/Os : $T(n) \leq O(n/3) + T(2n/3) = O(n/3)$

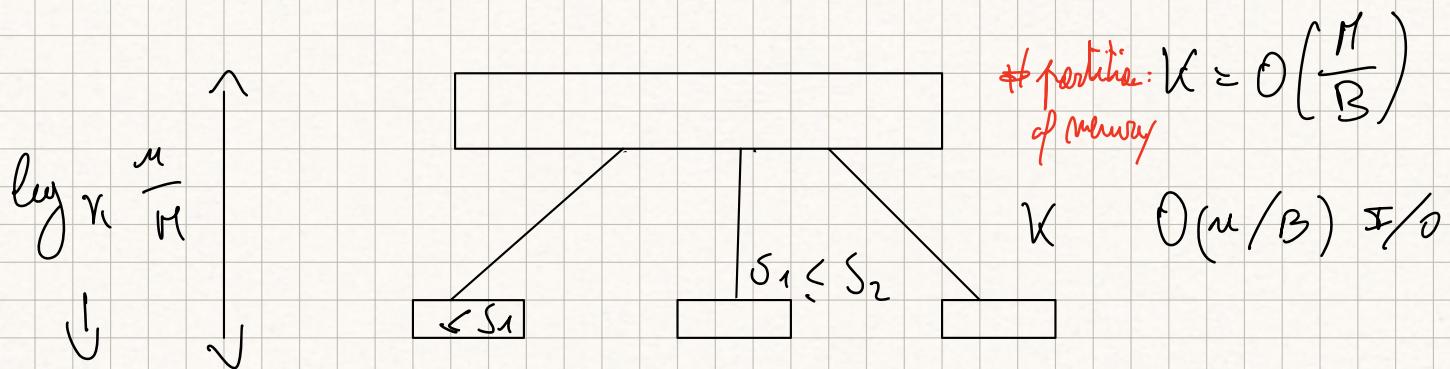
increasing n the complexity can not reduce !

Multi-Array Quicksort:



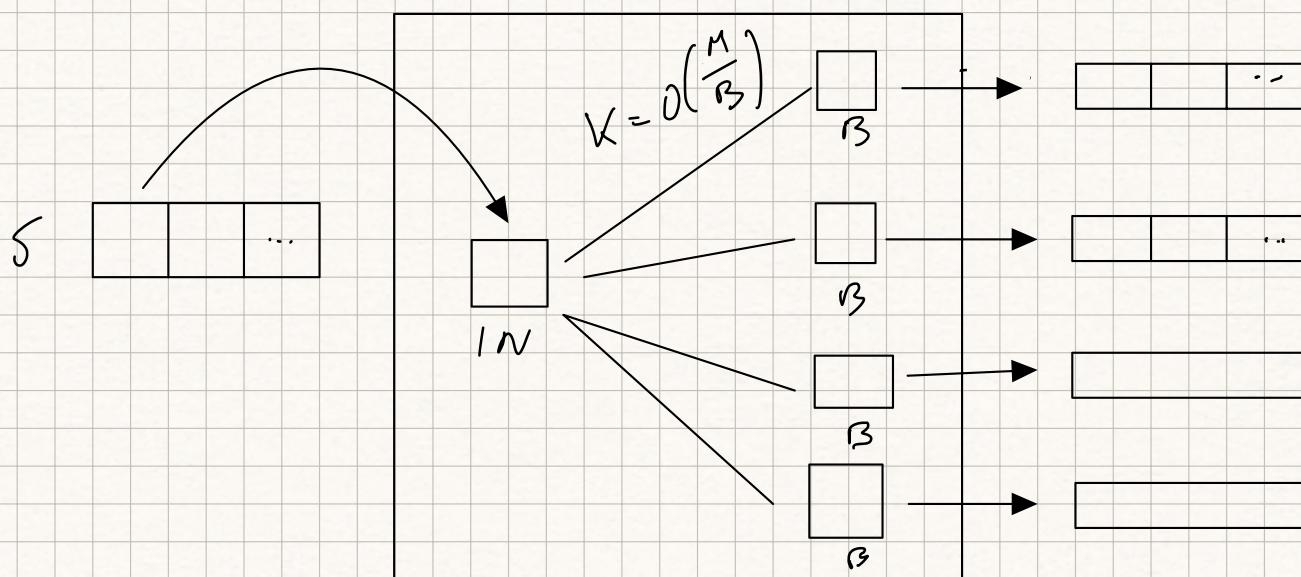
① more pivots s_1, s_2, \dots, s_{k-1}
 $s_1 = -\infty \quad < \quad s_2 \quad < \quad \dots \quad < \quad s_k = +\infty$

② $B_i = \{x \in S : s_{i-1} < x \leq s_i\}$ Buckets $|B_i| = O\left(\frac{n}{k}\right)$



when the tree is balanced B B B
 $\# \text{ of partitioning phase} = O\left(\log \frac{M}{B}\right)$

$$K \cdot B \leq M \rightarrow K \leq \frac{M}{B}$$

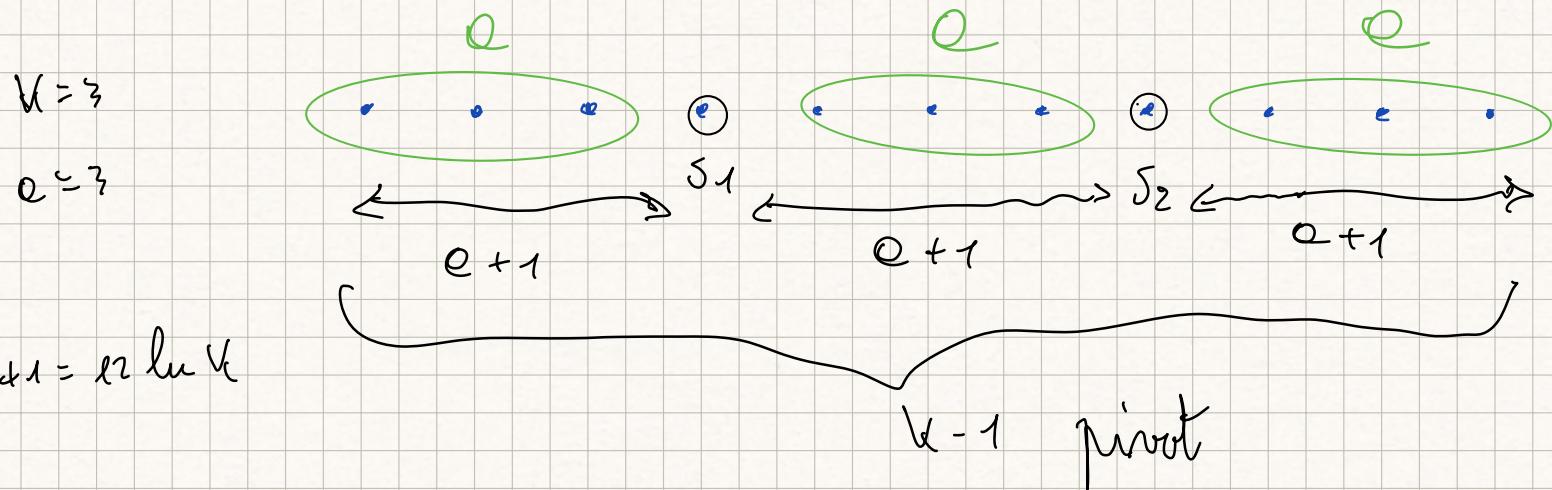


$$\Theta\left(\frac{n}{k} \cdot \log_{\frac{n}{k}}\left(\frac{n}{m}\right)\right) \Sigma/O_S \text{ OPT/MAL}$$

iff partition are balanced ($\sim \frac{n}{k}$ items each)

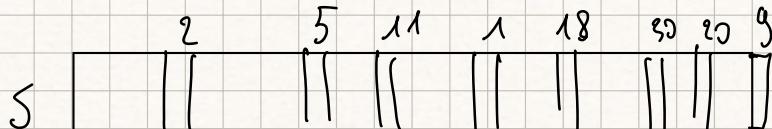
Oversampling (Multiway Quicksort):

- 1) Pick k $(e+1) \cdot k - 1 = (e+1)/(k-1) + e$ samples from $S \Rightarrow A$
- 2) Sort A $O = O(\log k)$
- 3) Pick pivot: $s_i = A[(e+1)i]$ $i = 1, \dots, k-1$



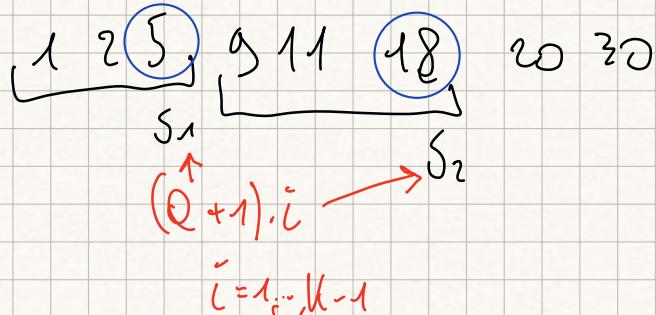
$$B_i = \{x \in S : s_{i-1} < x \leq s_i\} \quad s_0 = -\infty \quad s_k = +\infty$$

$$\text{Ex: } K=3 \\ e=2$$



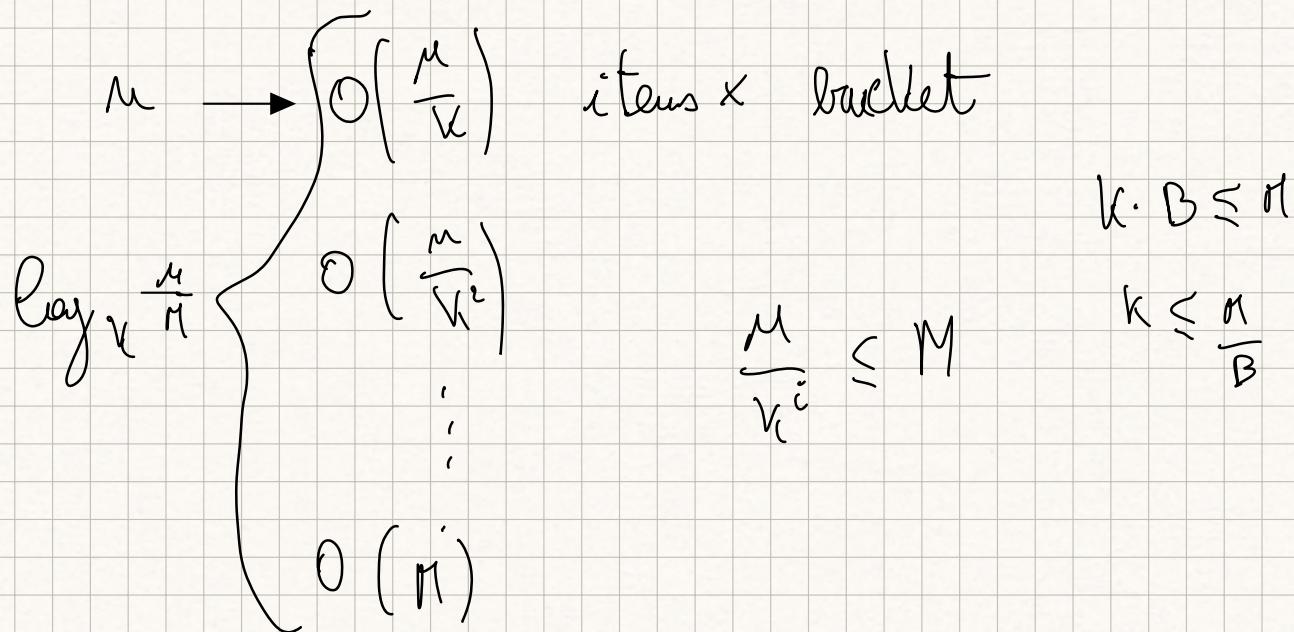
8 samples from S

$$(e+1)(K-1) + e = (e+1) + e = 8$$



Th: Let $K \geq 2$ and $\alpha + 1 = 12 \ln K$, a sample of size $(\alpha+1)K - 1$ suffices to ensure that all buckets get $\leq \frac{4n}{K}$ items with prob. $\geq \frac{1}{2}$

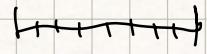
$$|B_i| \leq \frac{4n}{K} \quad \forall i=1, \dots, K$$



Dim: $P\left(\exists B_i : |B_i| > \frac{4n}{K}\right) \leq 1/2$

$$S_{\text{start}} = \left[\frac{2n}{K}, \frac{2n}{K}, \frac{2n}{K}, \dots \right]$$

$$t_1, t_2, t_3, \dots, t_{\frac{K}{2}}$$



$$S_{i-1} \quad B_i \quad S_i$$

$\frac{K}{2}$ blocks \rightarrow size of each block $= \frac{2n}{K}$

samples in $B_i = \varnothing$

$$\leq P\left(\exists t_j : t_j \text{ contains } < \alpha+1 \text{ samples}\right)$$

$$\leq \left(\frac{K}{2}\right) \cdot P(t_j : t_j \text{ contains } < \alpha+1 \text{ samples})$$

at the end

$$\leq \frac{K}{2} \cdot \frac{1}{K} = \frac{1}{2}$$

probability of 1 object

number of segment

① $P(\text{Q sample occurs in } t_3) = \frac{2^m/k}{m} = \frac{2}{k}$

② Fix $X = \# \text{ samples in } t_3 \rightarrow E[X] = \frac{2}{k} [(a+1) \cdot k - 1] = 2(a+1) - \frac{2}{k} \geq 2(a+1) - 1 \geq \frac{3}{2}(a+1) \rightarrow a+1 \leq \frac{2}{3} E[X]$

③ $P(X < a+1) \leq P(X < (1 - \frac{1}{3}) E[X]) \leq$

Chernoff bounds $P(X < (1 - \delta) E[X]) \leq e^{-\frac{\delta^2}{2} E[X]}$

$$\leq e^{-\frac{1}{3} \cdot \frac{1}{2} \cdot E[X]} \leq e^{-\frac{2}{3} \cdot \frac{1}{18} (a+1)} = e^{-\frac{1}{9} (a+1)} = e^{-\frac{1}{9} \cancel{12} \ln k}$$

$$= e^{-\ln k} = 1/k$$

Plugging $1/k$ in $P(\exists B_i : |B_i| \geq m/k) \leq (k/2) \cdot (1/k) = 1/2$ ■

- Given a sequence of S items, $S = (i_1, i_2, \dots, i_n)$ and a positive integer $m \leq n$, the goal is to select m items from S uniformly at random ($P = m/n$)

- $\text{Rand}(a, b) =$ returns a random number in $[a, b]$
- Assume that $m \leq \frac{n}{2}$

Dictionary = \emptyset

while $|D| < m$ do

$P = \text{rand}(1, m+1)$

if $P \notin D$ then insert (P, D)

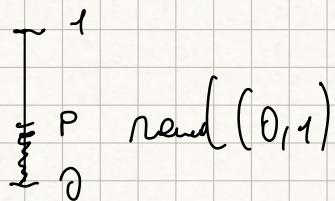
$$P(\text{collision}) = \frac{m}{m} \leq \frac{m}{n} \leq \frac{1}{2}$$

(take the same element)

- Assume that n is known

Take A with prob P

take B with prob $(1-P)$



$$S = i_1, i_2, i_3, i_4, \dots, i_n$$

$$P = \frac{1}{n}, \frac{1}{n-1}, \frac{1}{n-2}, \dots, \frac{1}{1}$$

$$m = 1$$

$$P = 1/m$$

$$P(\text{Pick } i_j) = \frac{1}{m-j+1} \quad P(\text{Pick } i_1, i_2, \dots, i_{j-1}) = \frac{1}{m}$$

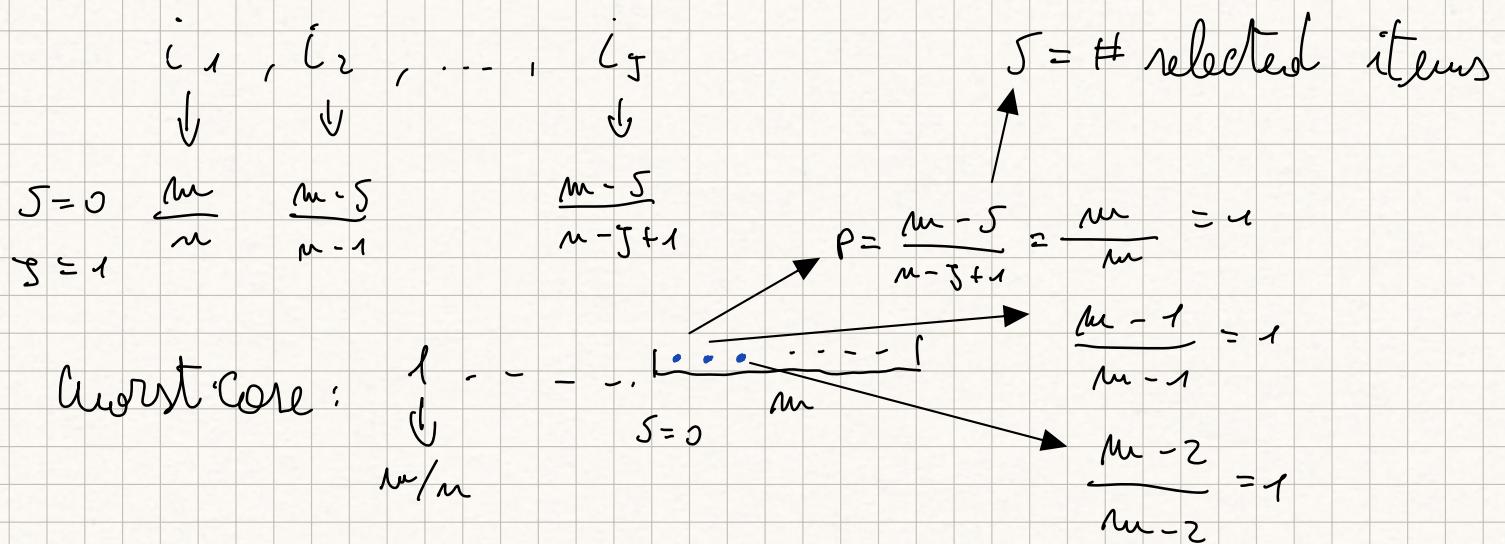
$$P(\text{Pick } i_1) = \frac{1}{m}$$

$P(\text{pick } i_j \text{ given that } i_1, \dots, i_{j-1} \text{ have not been picked}) =$

$$P(\text{pick } i_j) \cdot P(\text{not pick } i_1, \dots, i_{j-1}) = \frac{1}{m-j+1} \cdot \left(1 - \frac{j-1}{m}\right) = \frac{1}{m-j+1} \cdot \frac{m-j+1}{m} =$$

$$= \frac{1}{m}$$

$\forall m \geq 1$



Ex: $m = 3$ $i_1, i_2, i_3, i_4, i_5, i_6, i_7$ I will surely pick this 3 element

$m = 7$

$\frac{1}{7}, \frac{2}{7}, \frac{3}{7}, \frac{4}{7}, \frac{5}{7}, \frac{6}{7}, \frac{7}{7}$

$$P(\text{pick } i_j) = \frac{1}{m-g+1} \quad g=1, \dots, m$$

$\xrightarrow{i_1, i_2, \dots, i_m}$

↑
known

1) select the indexes $\{3, 7, 13\}$

2) pick items $\min\left\{m, \frac{m}{B}\right\}$

if m unknown: $P(\text{pick } i_g) = \frac{1}{m} \quad \forall g$ to be guaranteed

$0, 1 \quad 0, 7 \quad 0, 55 \quad 0, 05 \quad 0, 9 \quad 0, 09$
 $i_1, i_2, i_3, i_4, i_5, i_6$

random number $0 < r_g < 1$

We define $\forall i_g \Rightarrow \langle r_{g,m}, g \rangle \text{ Rand}(0,1)$ to get probability.

H : keep including the m priors of minimum priority (r_g)

Q: $P(\text{pick } i_g)$

$0, 1 \quad 0, 7 \quad 0, 55 \quad 0, 05 \quad 0, 9 \quad 0, 09$
 $i_1, i_2, i_3, i_4, i_5, i_6 \dots$

$H: \langle 0, 1, 1 \rangle$

$\langle 0, 7, 2 \rangle \leftarrow \langle 0, 55, 3 \rangle$ gets included because $0, 55 < 0, 7$ $\langle 0, 55, 3 \rangle \leftarrow$

Reservoir sampling:

$R = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 2 & \dots & \dots & m \\ \hline \end{array}$

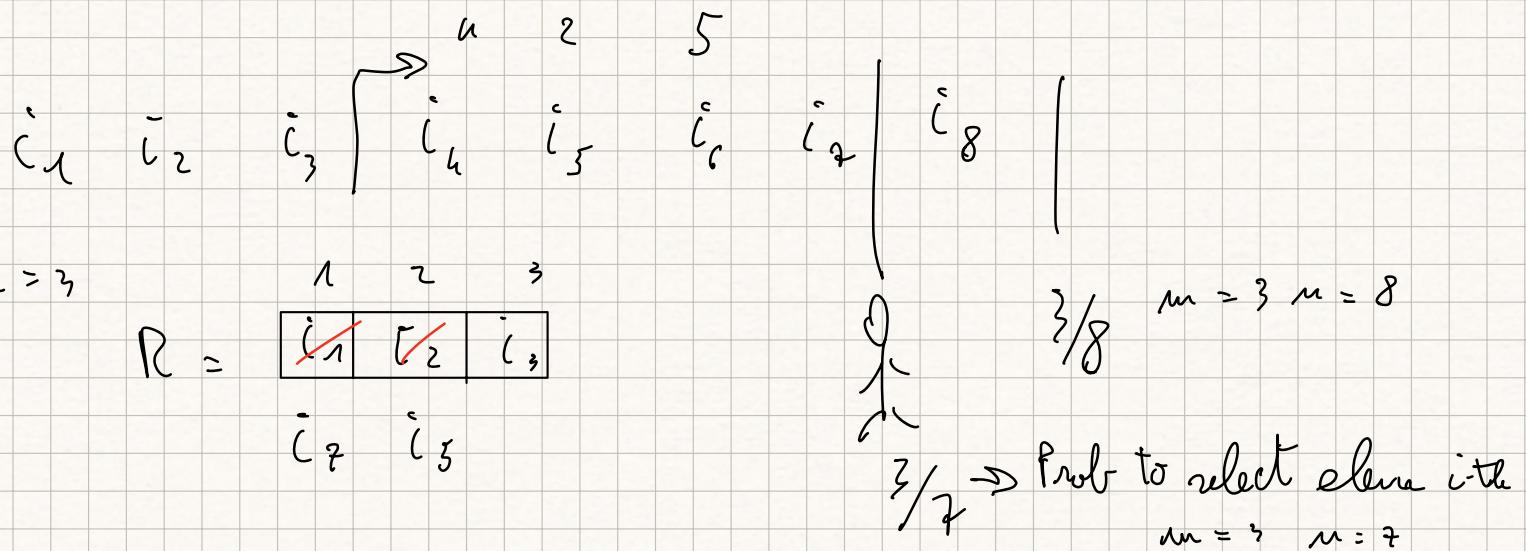
$\langle 0, 05, 4 \rangle$

1) initialize $R = i_1, i_2, \dots, i_m$

2) For $j = m+1$

$h = \text{rand}(1, j) \in N$

if ($h \leq m$) $R[h] = i_j$



Th: $\forall j \leq n : P(\text{pick } i_j) = \frac{m}{n}$ (UNIFORMLY AT RANDOM)

Dim:

bare core : $m = m$ $P(\text{pick } i_j) = 1$ $j = 1, \dots, m$ trivial

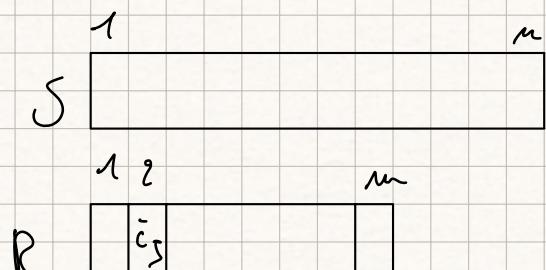
inductive core : property of uniform random sampling
is true for $m - 1$:

$$P(\text{pick } i_j) = \frac{m}{m-1} \quad j = 1, 2, \dots, m-1$$

Prove that this is true for m items $P(\text{pick } i_j) = \frac{m}{m}$
 $j = 1, \dots, m$

$$P(\text{pick } i_m) = \frac{m}{m} \text{ by design}$$

$$P(\text{pick } i_j) \text{ where } j = 1, 2, \dots, m-1$$



$P(i_j \text{ is in } R \text{ after } m-1 \text{ steps AND } ((i_m \text{ is not picked}) \text{ OR } (i_m \text{ is picked but stored not in position of } i_j)))$

initial hyp: $P(i_j \text{ picked}) = \frac{m}{m-1}$ for $j=1, 2, \dots, m-1$

$P(i_g \text{ is in } R \text{ after } m-1 \text{ steps}) \cdot P(\text{not picked}) \text{ or}$
 $(\text{is picked} \wedge \text{not stored where is } i_g)$

$$\frac{m}{m-1} \cdot \left[\left(1 - \frac{m}{m} \right) + \left[\frac{m}{m} \cdot \frac{m-1}{m} \right] \right] =$$

$$= \frac{m}{m-1} \cdot \left(\frac{m-m}{m} + \frac{m-1}{m} \right) = \frac{m}{m-1} \cdot \left(\frac{m-1}{m} \right) = \frac{m}{m}$$

PROBLEM: we need to generate two random numbers.

e.g.: $m=2$ $n=8$ known \rightarrow sequence length

$1, 2, 3, n, 5, 6, 2, 8 \leftarrow j$
 $S = (e, d, c, b, a, f, g, h)$

$$P \leq \frac{m-5}{m-j+1}$$

$$P = (.5, .5, 0, 1, .5, 1, 0, 02, 1, 1)$$

$$R = \boxed{\quad} \quad \boxed{\quad}$$

$$j=0 \quad j=1 : \frac{m-5}{m-j+1} = \frac{2-0}{8-1+1} = \frac{2}{8} = 0,25 \times$$

$$j=0 \quad j=2 : \frac{2-0}{8-2+1} = \frac{2}{7} \approx \frac{1}{2} \times$$

$$S=0 \quad J=3 : \frac{2-0}{8-3+1} = \frac{2}{6} = \frac{1}{3} \sim 0,1 \quad \checkmark$$

$$S=1 \quad S=h : \frac{2-1}{8-h+1}$$

$$S = (0, b, c, d, e, f, g, h, i, \dots) \quad m=3$$

$$d_h = (\dots \ 2 \ 4 \ 1 \ 2 \ 3 \ 1) \quad \begin{matrix} \uparrow \\ \text{status } R \end{matrix}$$

$$R = \boxed{0 \ b \ e} \quad \begin{matrix} 1 & 2 & 3 \end{matrix}$$

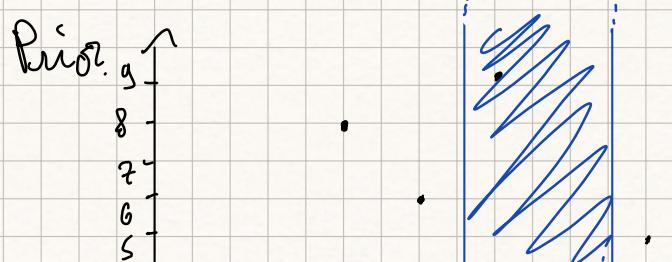
~~f d h~~

~~i g~~

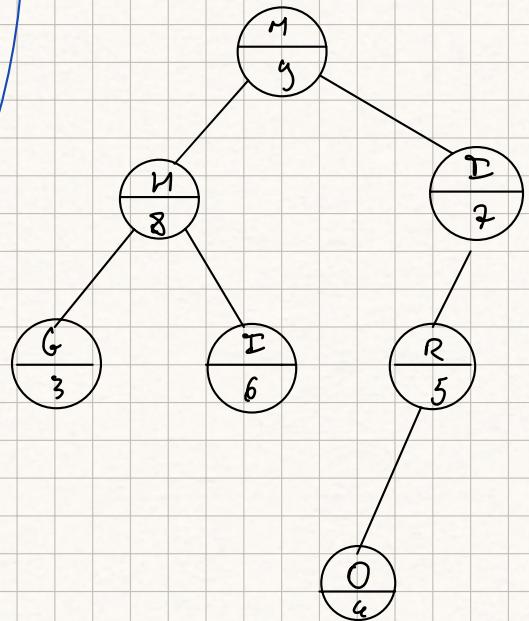
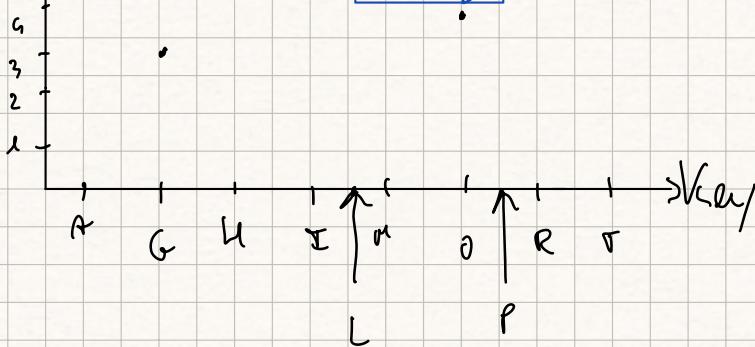
heap
T R E A P
Free

min
max

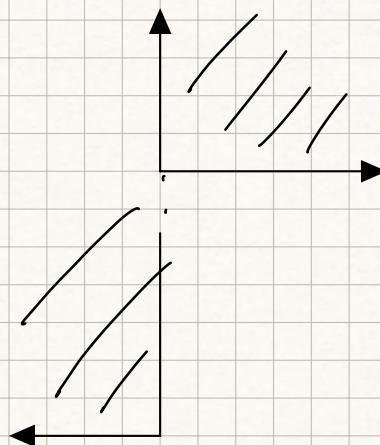
Key = letter
Priority = integer



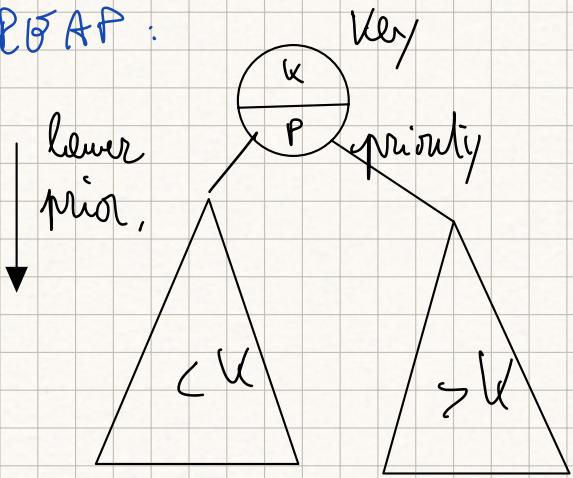
return all pairs in
 $[l, r] \times [s, +\infty)$



3 sided range query



TPOFAP :



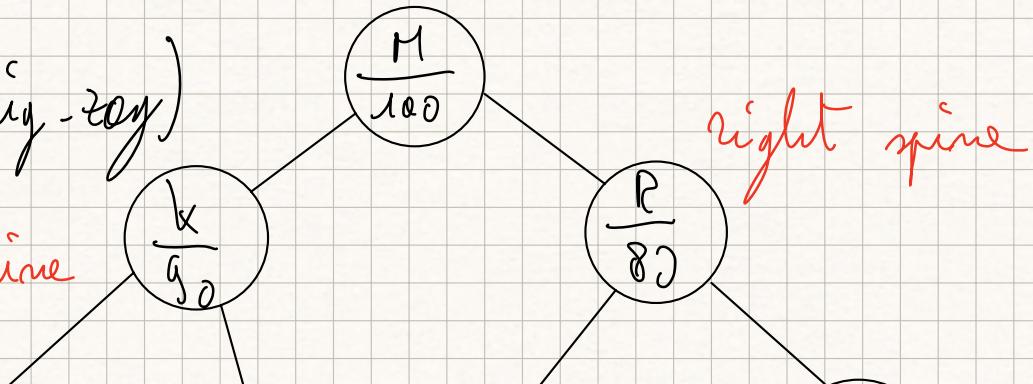
- Binary search tree on Keys
- Keep (max) on priorities

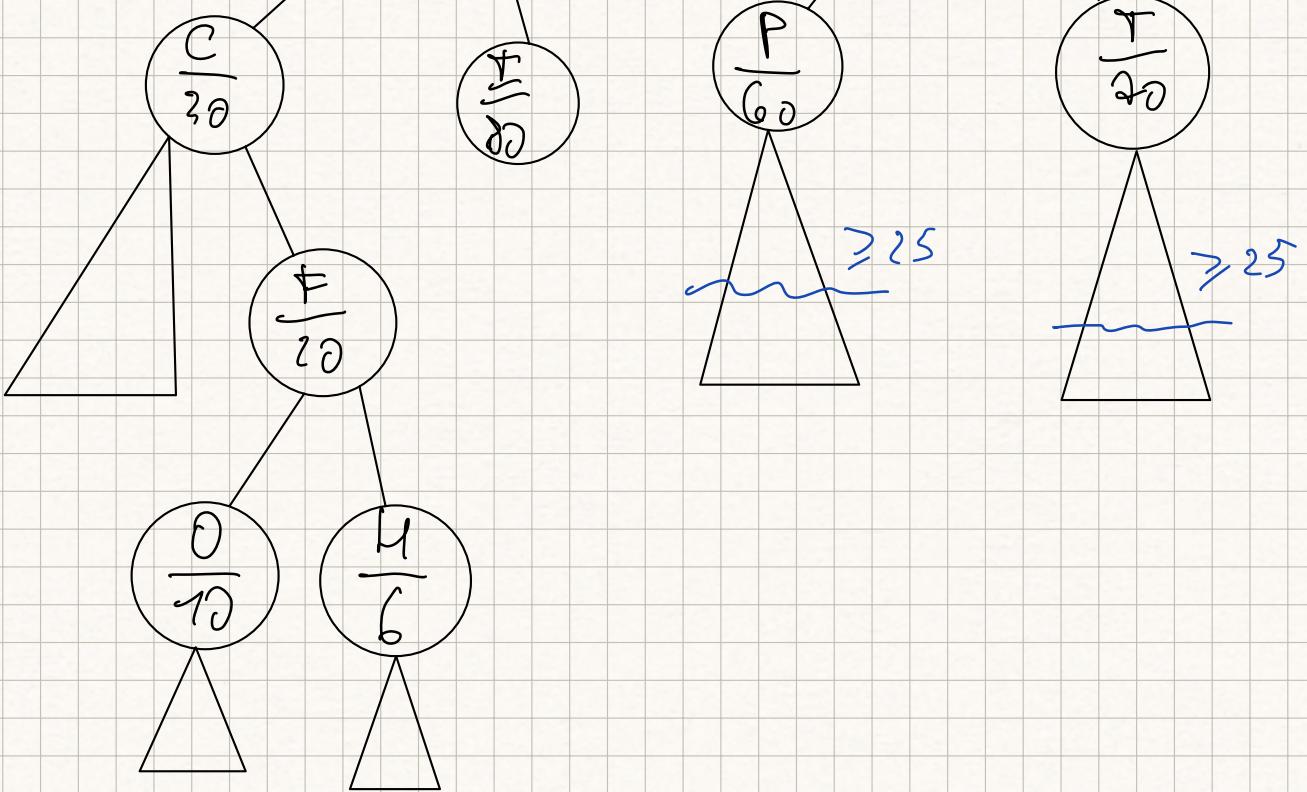
$$[G, T] \times [2S, +\infty)$$

x - axis y - axis

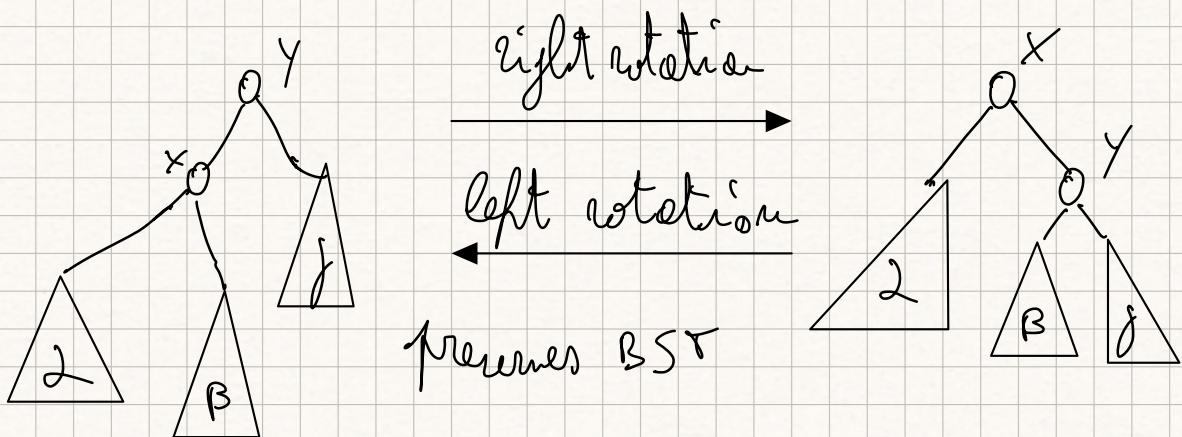
$O(h)$ = traversal (zig-zag)

left spine





$O(h) +$ ↗ # returned answers
 ↗ # visited but stopped nodes
 ↗ 2. ↗ # returned answers



$$x \leq y$$

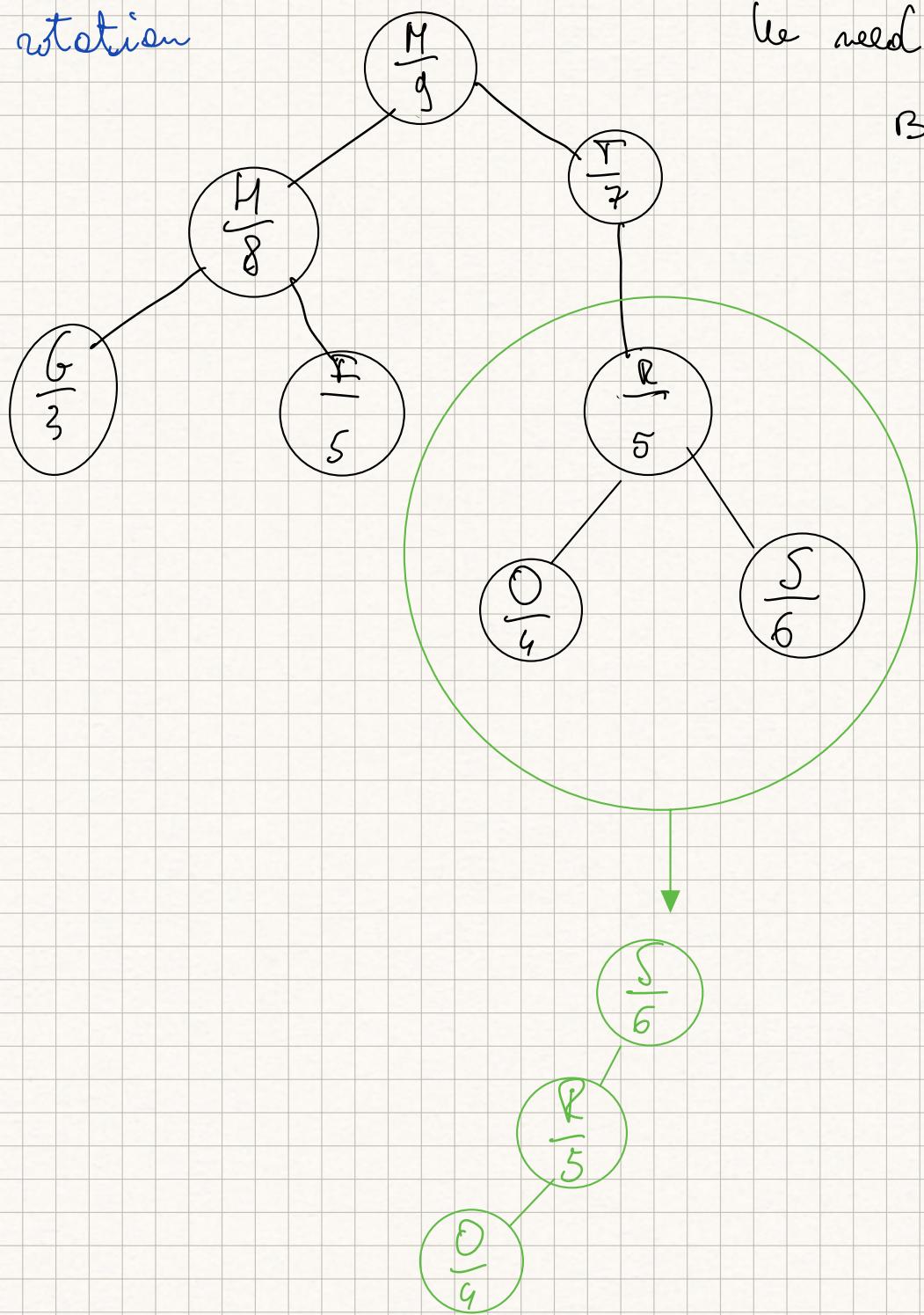
presumes x
denotes y

Search(V) search in a BST (binary search tree)

Insert(V, P): $O(h)$ $\frac{V}{P}$ $\rightarrow T$

- 1) Search
- 2) rotation

We need to preserve the BST & keep



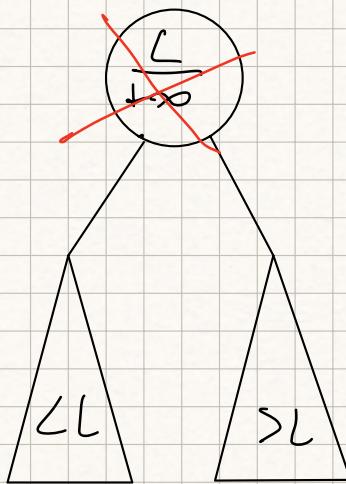
$\text{Delete}(v) : O(h)$

- 1) search
- 2) set $-\infty$ priority to make to del
- 3) rotate until it has no child.

$\text{Split}(v) : O(h)$

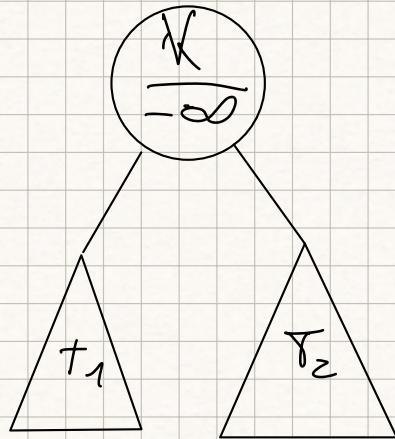
$\text{split}(L) \rightarrow T_{\leq L}, T_{> L}$

- Insert $(L, +\infty)$



$\text{Merge } (+_1, T_2) : O(h)$

Keeps $|+_1| < \text{Keys}(T_2)$



$$T_1 < K < T_2$$

rotate for the root $\frac{K}{\infty}$ until it has no child

$\text{Split}(v) = \langle T_1, T_2 \rangle$

$O(h)$

Insert a dummy node with v as key and $+\infty$ priority and rotate it to the radix and then delete it

3-side query = $Q : [a, b] \times [c, +\infty)$

$\text{Split}(e)$

keys priority

$(-\infty, 0]$

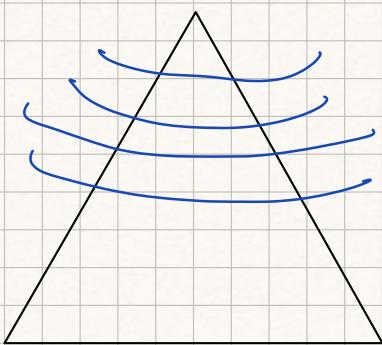
$(0, +\infty)$

$\text{Split}(d, T_{>0})$

$[e, f]$

$(f, +\infty)$

top-bottom approach



ex: $[c, 0] \times [6, +\infty)$ on first sweep

Given a dictionary D , build a BST balanced

- AVL-tree, 2-3 trees, Red-Black trees, ... h: $O(\log n)$

Th: Build a tree in which keys $\in D$, priorities are

random numbers. height = $O(\log n)$ on expectation

Def:

$$A_{ik}^i = \begin{cases} 1 & \text{if key } x_i \text{ is a proper ancestor of key } x_k \\ 0 & \text{otherwise.} \end{cases}$$

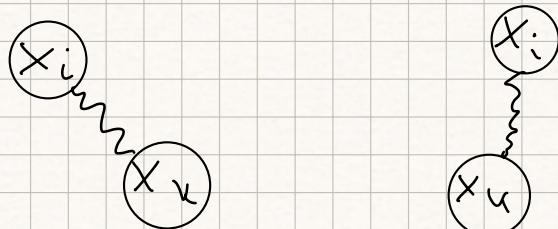
random indicator variable

$$D = \{x_1, x_2, \dots, x_n\}$$

$$\text{depth}(x_k) = \sum_{i=1}^n A_{ik}^i \quad \# \text{ of ancestor of } x_k$$

$$\mathbb{E}[\text{depth}(x_k)] = E\left[\sum_{i=1}^n A_{ik}^i\right] = \sum_{i=1}^n E[A_{ik}^i] =$$

$$= \sum_{i=1}^n P(x_i \text{ is a proper ancestor of } x_k) = \sum_{i=1}^n 1 \cdot P(A_{ik}^i = 1)$$



max heap: $\text{prior}(x_i) > \text{prior}(x_k)$

min heap in
the nodes

1) $X(i, k) = \{x_i, x_{i+1}, \dots, x_k\} \quad i < k \quad \# \text{items} = k - i + 1$

$x_1 \dots x_{i-1}$

$\rightarrow x_{i+1} \dots x_m$

2) $X(k, i) = \{x_k, x_{k+1}, \dots, x_i \mid i > k \quad \# \text{items} = i - k + 1\}$

$$\lceil \text{depth}(x_n) \rceil = \sum_{i=1}^n P\left(\frac{x_1}{x_n}\right) = \sum_{i < k}^{k-1} \frac{1}{k-i+1} + \sum_{i \geq k}^n \frac{1}{i-k+1} <$$

$i=1 \rightarrow 1/k \quad | \quad i=k+1 \rightarrow 1/2$

$i=2 \rightarrow 1/k-1 \quad | \quad i=k+2 \rightarrow 1/3$

$\vdots \quad \vdots \quad | \quad \vdots$

$i=k-1 \rightarrow 1/2 \quad | \quad i=n \rightarrow \frac{1}{n-k+1}$

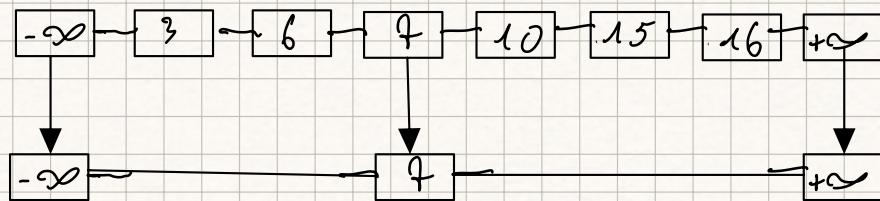
$\ln n \geq \ln k$

$$< \ln k + \ln(n - k + 1) - 2 < 2 \ln n - 2 \rightarrow O(\log n)$$

Skip list

- Level DB
- RODS (SortedSet)
- Rocks DB
- llBore

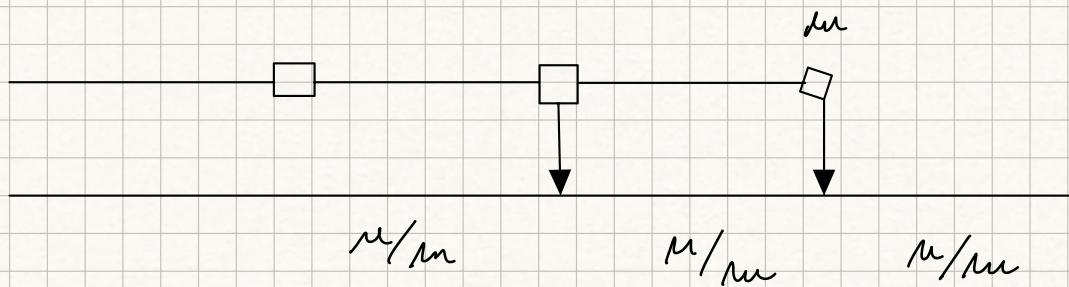
Do list



$$m = \# \text{ items} = |L_0|$$

$L_1 \text{ items} = m$

$L_0 \text{ items} = m$



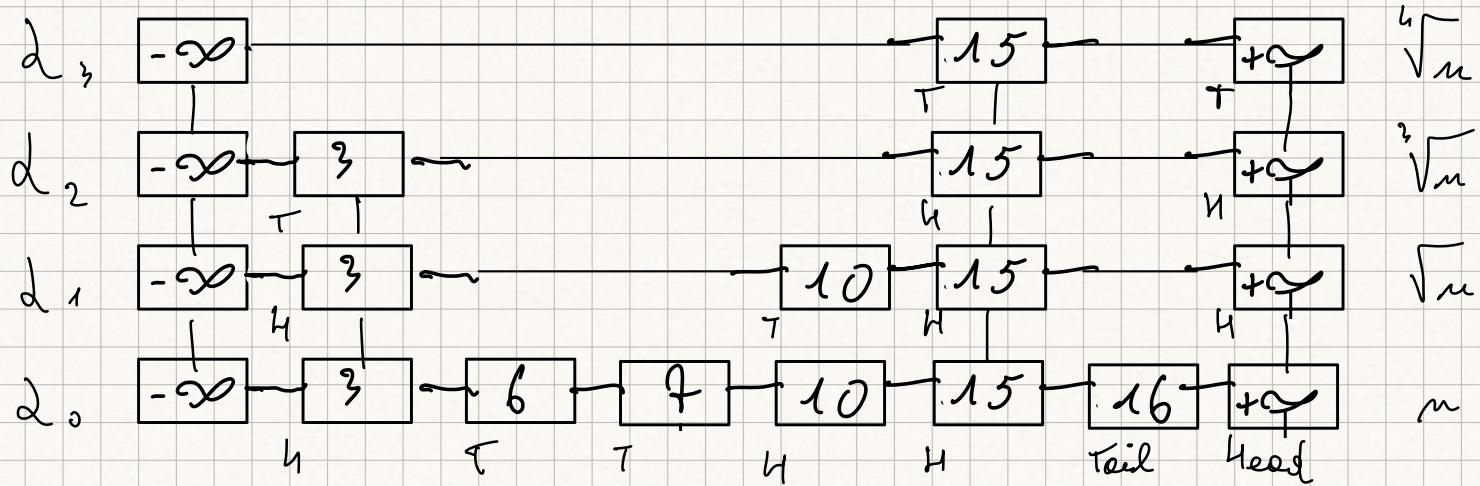
cost of a search $O(m + \frac{r}{m})$

'if' $m = \frac{n}{m}$ minimizing the search cost

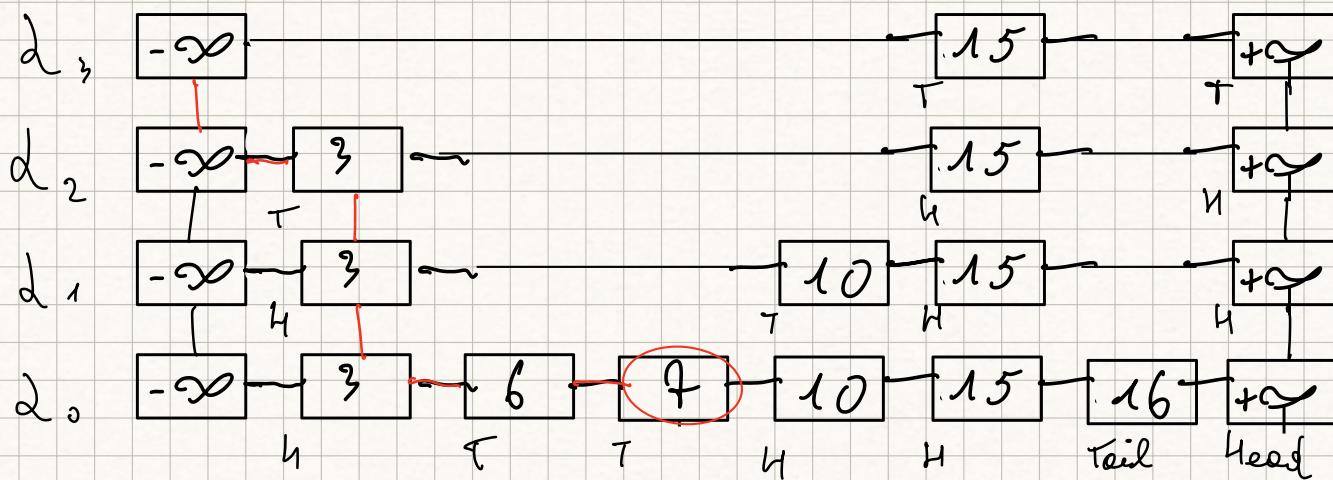
$$m^2 = n$$

$n = \sqrt{n}$ optimal choice

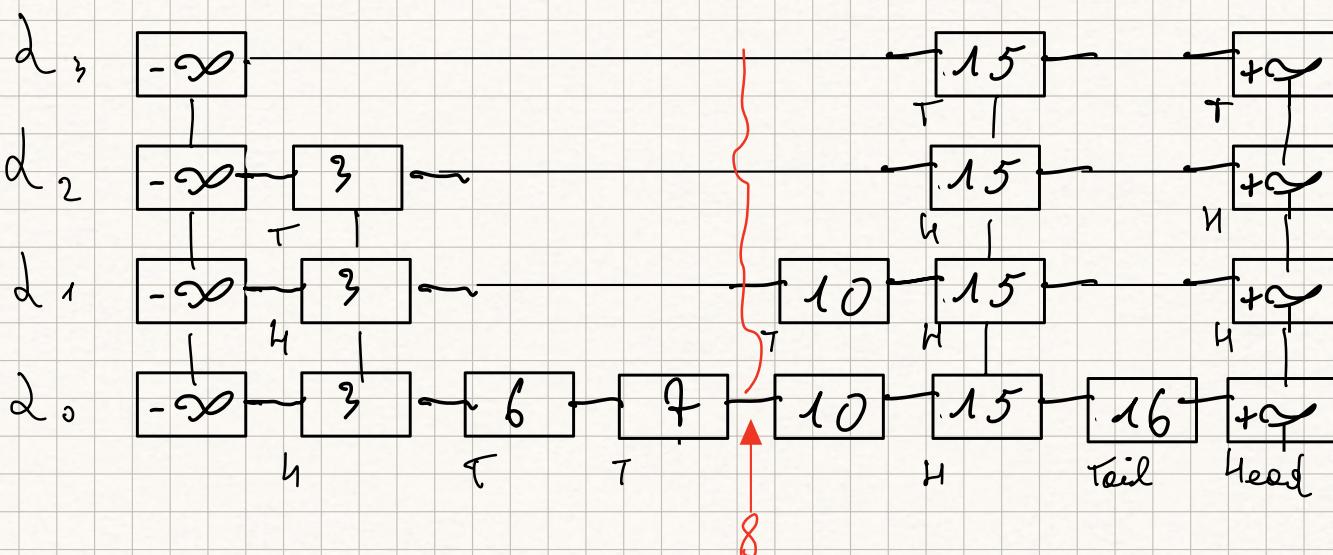
use 2 different approaches $E[\text{height}] = O(\log n)$



Search \neq in Skip list



Insertion of 8 in Skip list



Start from a coin and decide in which list put it

Time complexity of search = $\mathcal{O}(\# \text{vertical links traversed} + \# \text{horizontal links traversed})$

$$\text{level} = \max_{\text{key } k} L(k)$$

level for key k

$$P(L(k) \geq l) = \text{at least } l \text{ heads} = \left(\frac{1}{2}\right)^l$$

$$P(L \geq l) = P(\max_k L(k) \geq l) = P(\max\{L(1), L(2), \dots, L(n)\} \geq l)$$

union bound $\leq \sum_{k=1}^n P(L(k) \geq l) = \sum_{k=1}^n \left(\frac{1}{2}\right)^l = \frac{n}{2^l}$ Probability that height is greater or eq. than l

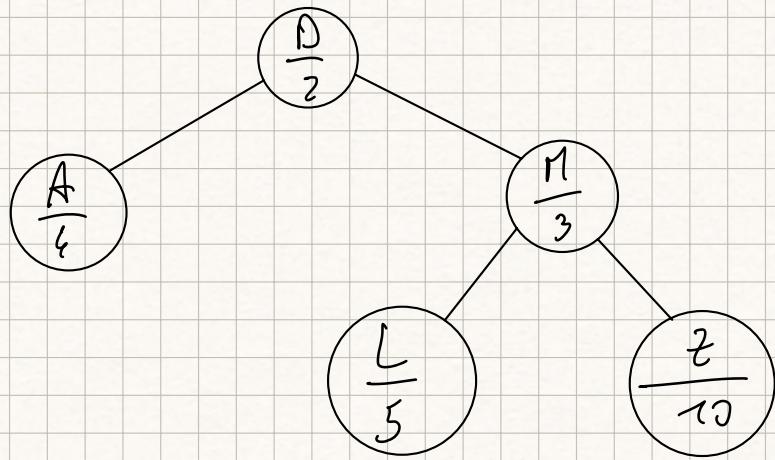
if $l \leq \log_2 n \rightarrow \frac{n}{2^l} \geq 1$ Trivial Prob. doesn't hold because > 1

if $l = c \cdot \log_2 n \rightarrow c > 1$

$$P(L \geq l) = \frac{n}{2^{c \cdot \log_2 n}} = \frac{n}{n^c} = \frac{1}{n^{c-1}} \quad c > 1 \quad \text{with high probability}$$

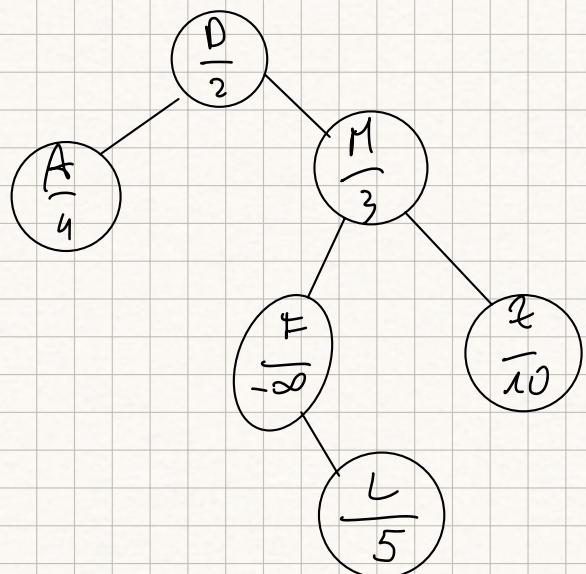
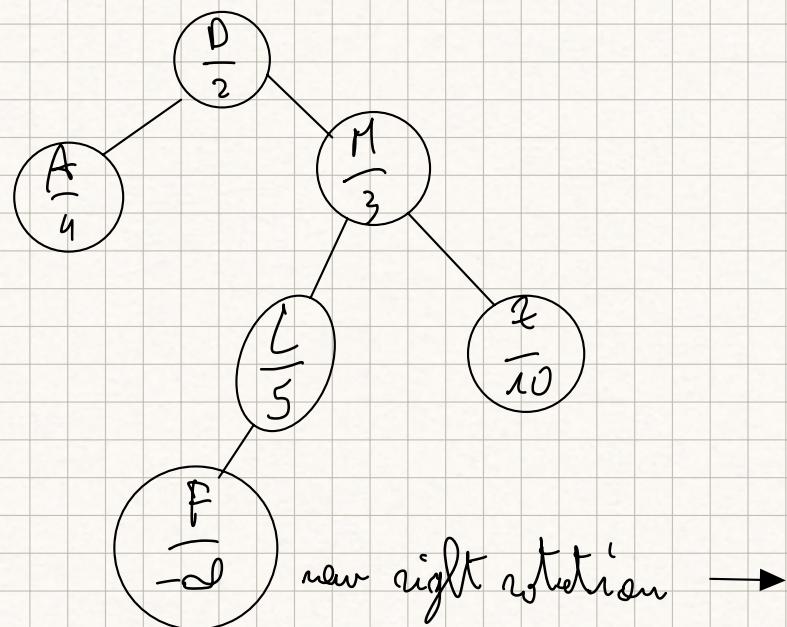
We conclude that with high probability, a skip list has $O(\log n)$ levels.

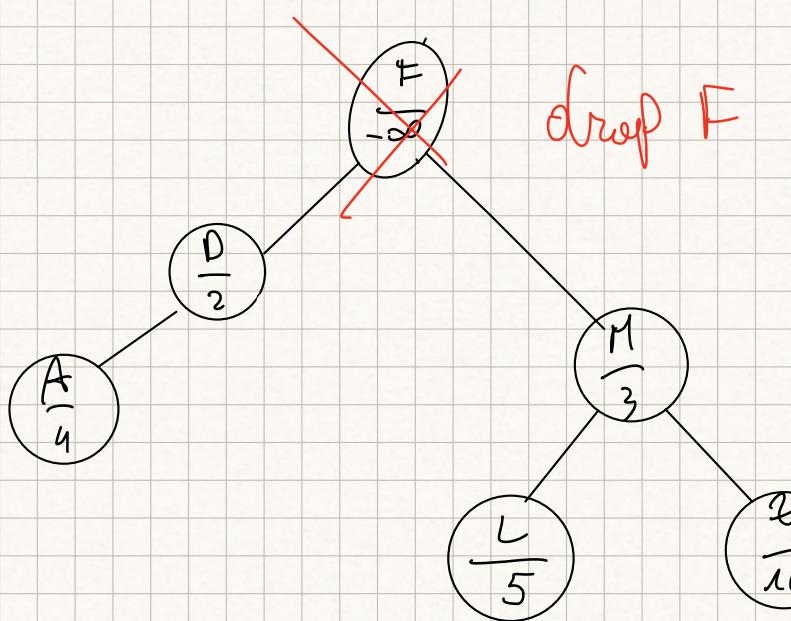
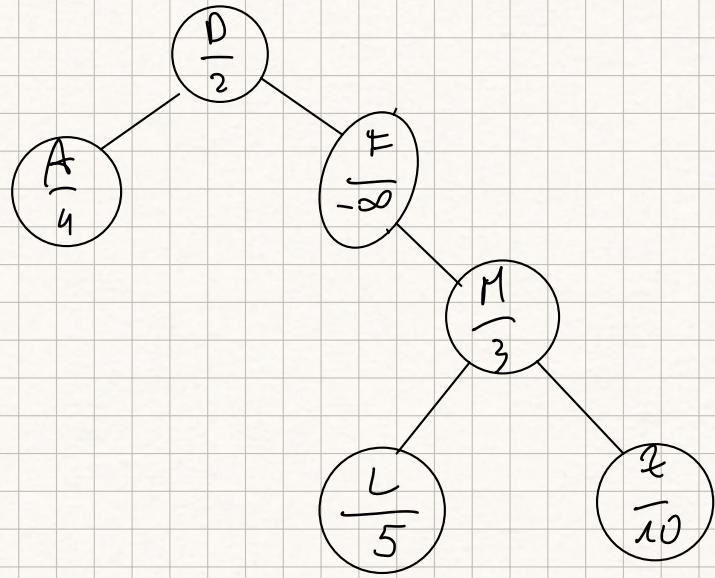
ds: TRAP(min)



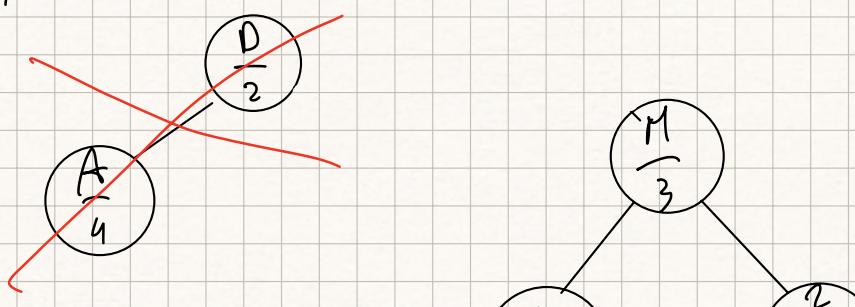
Split (F)

· Insert ($F, -\infty$)





3-node question $[F, v] \times [g, -\infty]$

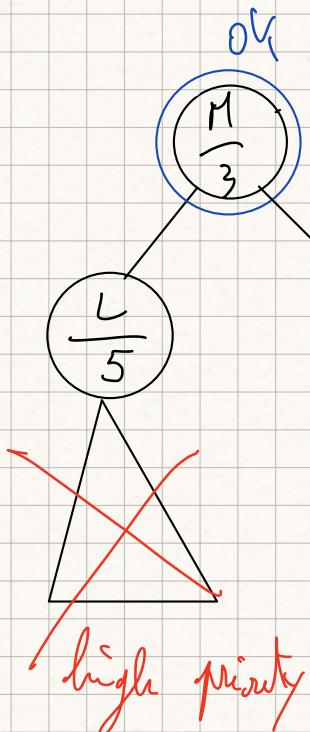
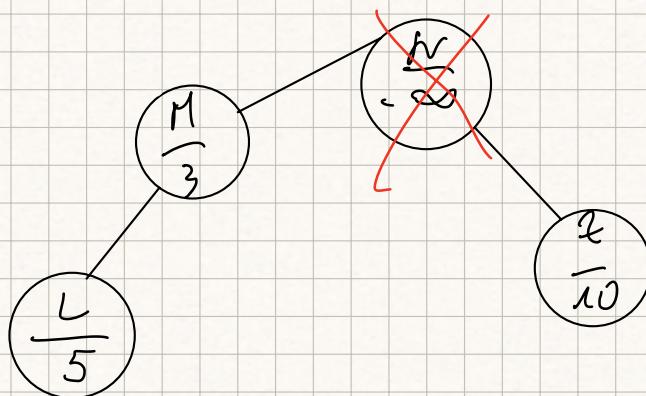
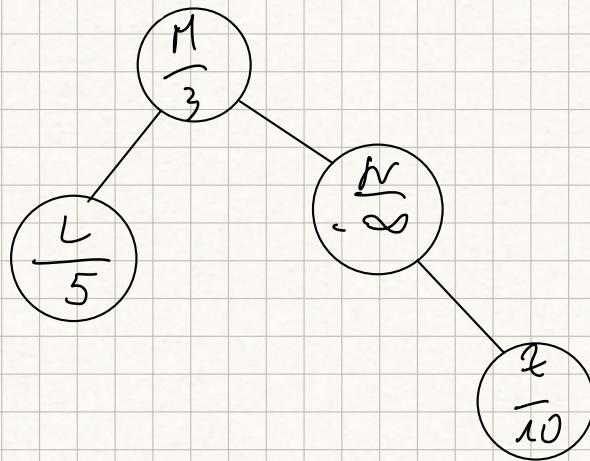
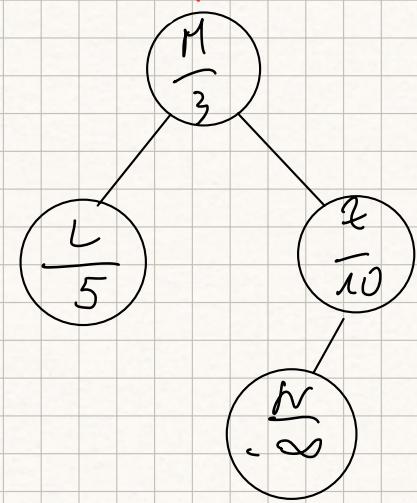


$$\frac{L}{5}$$

$$\frac{x}{10}$$

1 Split F (already done)

2 Drop left subtree ($i_{\text{ten}} < F$)



~~$\frac{x}{10}$~~
high priority

high priority

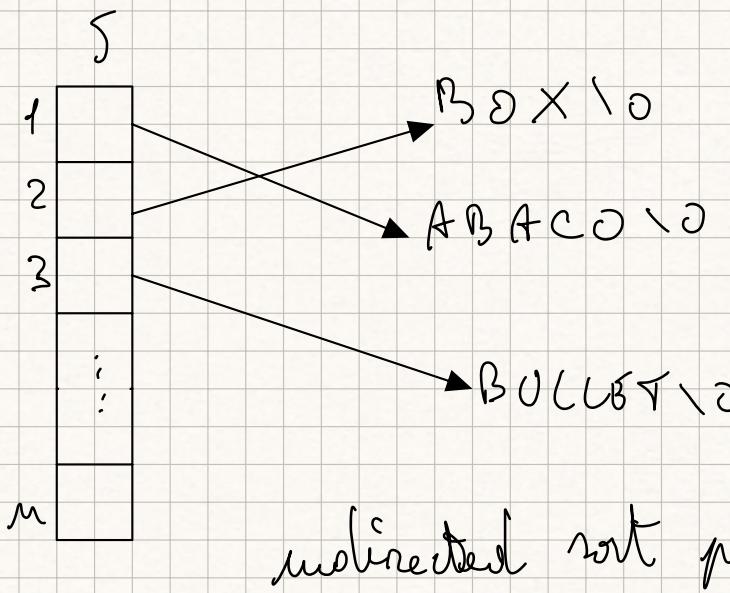
String sorting (RAM):

$S[1, n]$ = array of pointers to strings

N = Total length of strings

σ = alphabet size

$L = \frac{N}{m}$ average string length



qsort ($S, m, \text{char} \star, \text{cmp}$)

$O(m \log n)$ on expectation

$O((m \cdot \log n) \cdot l)$ time complexity
 avg string len = $\frac{N}{m} \rightarrow$ # of strings
 \rightarrow tot string len

$O(m \log n)$ cache misses / I/Os

distinguishing prefix of a string in a set S

$$S \left\{ \begin{array}{ll} S_1 = \text{abccr} & d_{S_1} = 4 \\ S_2 = \text{abate} & d_{S_2} = 4 \\ S_3 = \text{abt} & d_{S_3} = 3 \end{array} \right\} d = \sum_S d_S = 11$$

a) Only voter has to look at the distinguishing prefix of $s \in S$

$$\mathcal{O}(d)$$

b) Sorting S needs to sort by the first char of a string $\mathcal{O}(m \log n)$

Lower bound # comparisons to sort n strings \Rightarrow

$$= \mathcal{O}(0 + m \log n)$$

| m | $\log n$ |
|-------------|----------------------------|
| 0 0 - - . 0 | 0 0 $\downarrow \lim(0)$ |
| 0 0 - - . 0 | 0 1 $\downarrow \lim(1)$ |
| ⋮ | ⋮ |
| 0 0 - - . 0 | 0 111 $\downarrow \lim(n)$ |

$$\forall s : ds > m \Rightarrow O = \mathcal{O}(m \cdot n)$$

lower bound on that $S = \mathcal{O}(m \cdot n + m \log n)$ either empty

length of the strings we are comparing.

$$q\text{sort} = O((n \log n) \cdot (m + \log n))$$

$m = \text{constant} = O(1)$ lower bound = $\Omega(n \log n)$ $q\text{sort} = O(n \log^2 n)$

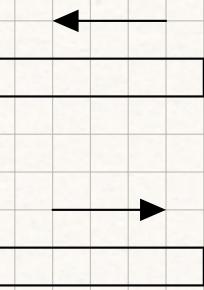
$m > \log n$ l. b. = $\Omega(m \cdot n)$ $q\text{sort} = O(m \cdot n \log n)$

$m \leq \log n$ lower bound = $\Omega(n \log n)$ $q\text{sort} = O(n \log^2 n)$

Q sort is not optimal. there is a factor n of $O(\log n)$ between lower bound and qsort

Radix-based sort:

- LSD : least-significant digit } first
- MSD : most-significant digit }



TRIV

If alg. will introduce a data structure,

α = distinguishing prefix

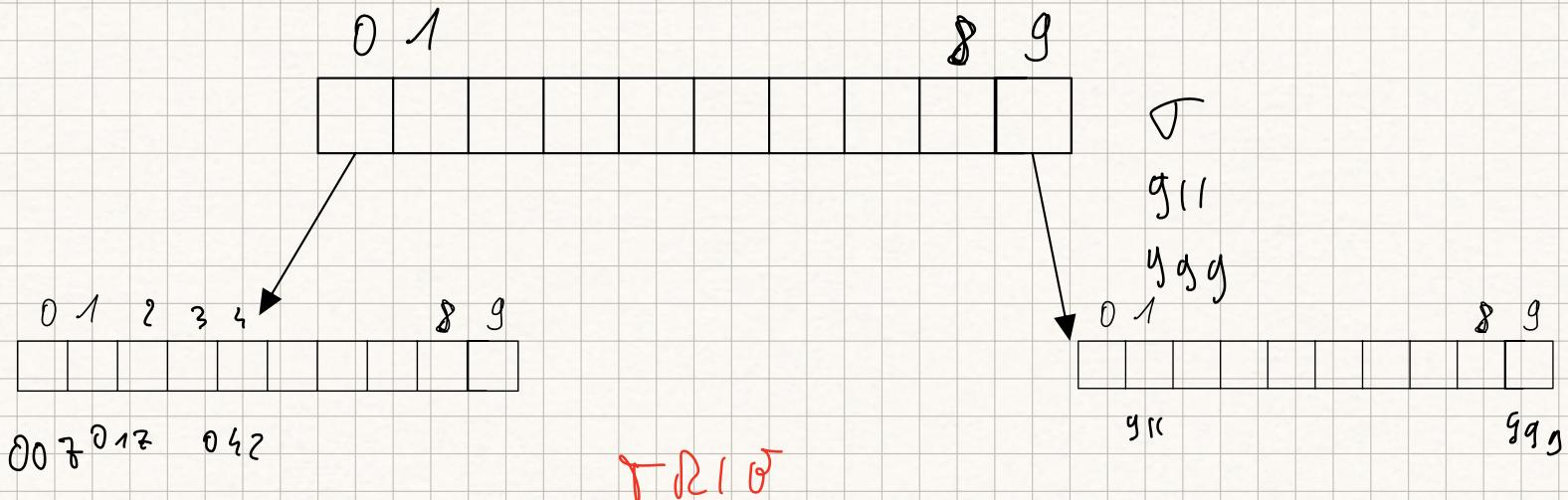
$n = \# \text{ strings}$

$r = |\Sigma|$ alphabet size

$N = \text{Total length of the strings}$

MSD

$$\Sigma = \{00^2, 017, 042, 911, 999\}$$



$$\text{size trie} = \# \text{ nodes} = O\left(\sum_s d_s\right) \in O(d)$$

leaves = $n = \# \text{ strings}$

Space = $O(\# \text{ nodes} \cdot \sigma) = O(d \cdot \tau)$ $\tau = 2^{64}$ a lot of ~~space~~ pointers

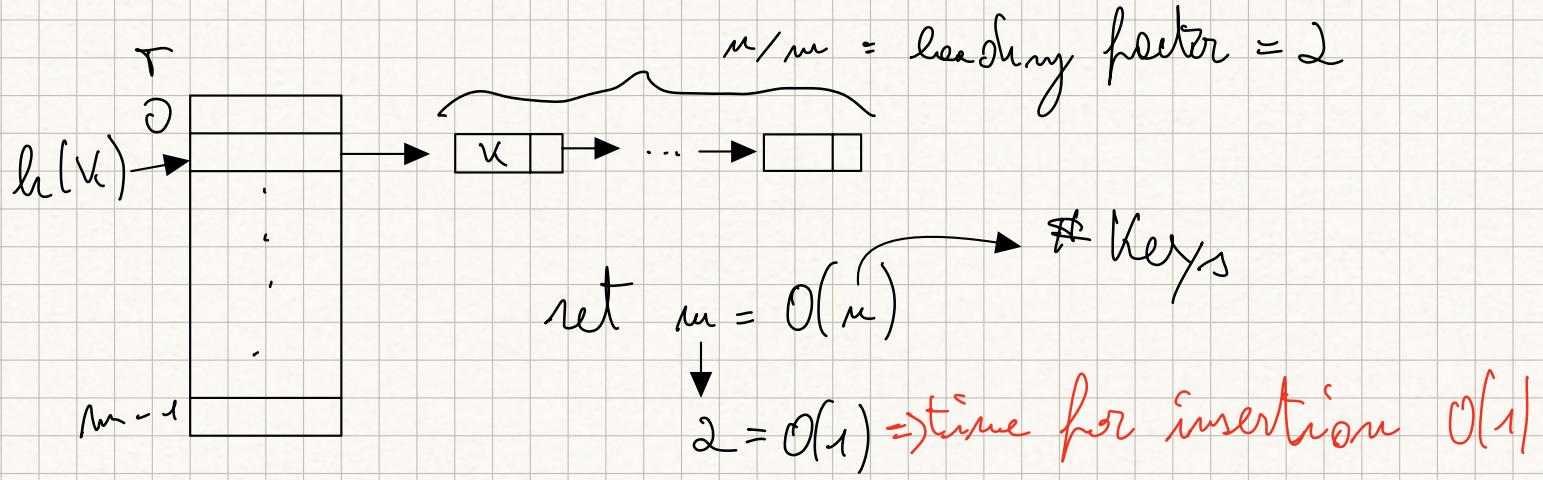
$$\text{Time construct trie} = O\left(\underbrace{d}_{\text{string dist.}} + \underbrace{d \cdot \tau}_{\text{node creation}} + \underbrace{(d \cdot \sigma + n)}_{\text{visit }}\right)$$

$O(d \cdot \tau)$

$d \approx n$

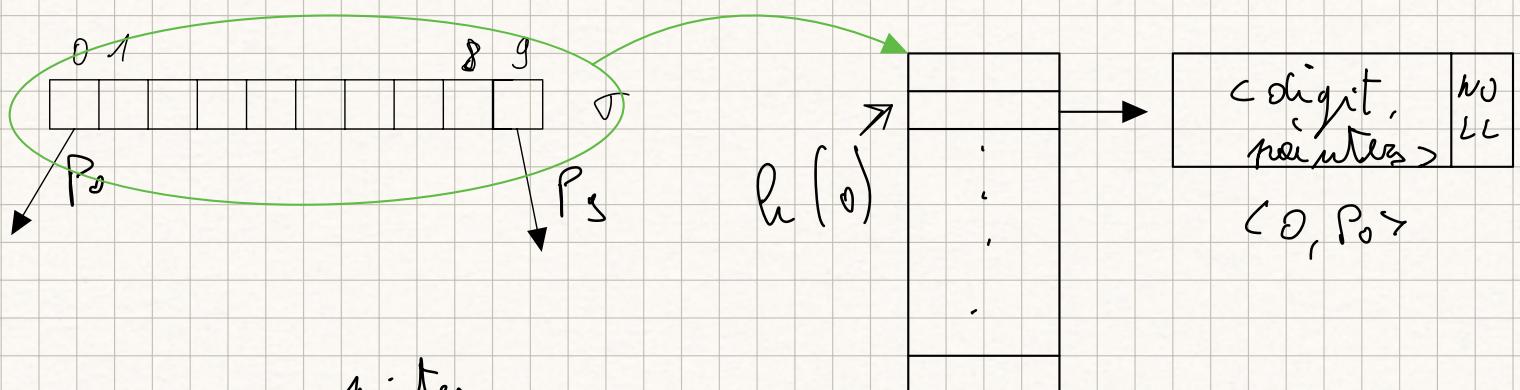
Prefix search on trie: $O(d_s)$

Hash Table (hashing with chaining)



Store k Keys taking $O(k)$ Time, and Space expected

Hash Table for Trie:



$$\text{Space} = O(m + \sum_m e_m) \leq m + d = O(d)$$

of edges having m
size of all hash tables $\leq m + d$

$$\begin{aligned}
 \text{Time} &= O\left(\underbrace{n + d}_{\substack{\text{construction} \\ \text{time}}} + \sum_m l_m \cdot \log r_m\right) \leq O(d \log \sigma) \\
 &\leq \sum_m l_m \cdot \log \sigma = (\log \sigma) \cdot \sum_m l_m = d \log \sigma
 \end{aligned}$$

$(d+m) \leq d > n$
 $\sum_m l_m = d$

Time to sort with an hash table + tree = $O((d+n) \cdot \log \sigma)$

LSD - first Radix Sort (use counting sort that is stable)

| | | | |
|-------|-------|-------|-------|
| 0 1 7 | 1 1 1 | 0 0 7 | 0 0 7 |
| 0 4 2 | 0 4 2 | 1 1 1 | 0 1 7 |
| 6 6 5 | 6 6 5 | 0 1 7 | 0 4 2 |
| 1 1 1 | 0 1 7 | 0 4 2 | 1 1 1 |
| 0 0 7 | 0 0 7 | 6 6 5 | 6 6 5 |

- L phases
- $O(n + \sigma)$ time per phase

$O(L \cdot (n + \sigma))$ time

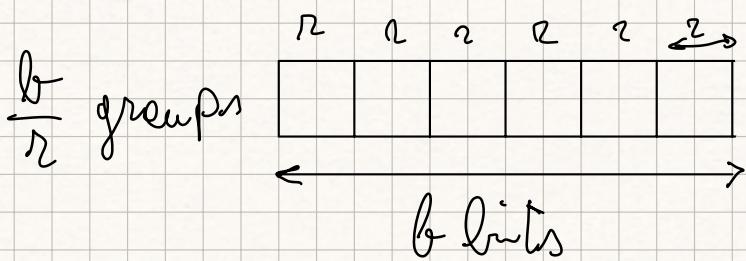
$$O(L \cdot n + n \cdot \sigma) = O(N) \text{ time}$$

for string len.

MSD $\begin{cases} O(d \log \sigma) \\ O(d) \text{ space} \end{cases}$

Assume that all keys are binary

- every phase cost: $O(n + 2^n)$ time



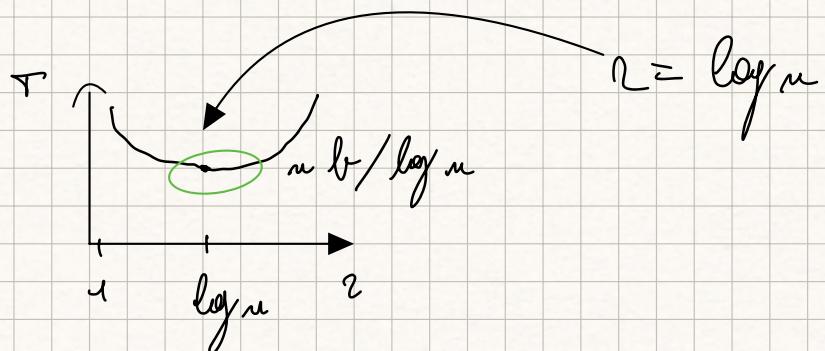
$2^n = \text{alphabet size}$

$$O\left(\frac{b}{n} \cdot (n + 2^n)\right) \text{ Time}$$

- # phases = $\frac{b}{n}$ = # of letters in a word

if: $2 \leq \log_2 n \rightarrow 2^n \leq n$ Time compl. = $O\left(\frac{b}{n} \cdot n\right)$

if: $n > \log_2 m \rightarrow 2^n > n$ Time compl. = $O\left(\frac{b \cdot 2^n}{n}\right)$



More time compl. $n = \log_2 m \rightarrow O\left(\frac{\log n}{\log_2 m}\right) \text{ Time} = O\left(\frac{N}{\log_2 m}\right) \text{ Time}$ ■

$$\frac{N}{\log_2 m} \sim \text{st. log } \sigma$$

MSD

Multi-Vary QuickSort: comp-based sorting for strings:
 $O(n \log n + d)$

$MVK_Q(R, i)$ ($i \geq 1$; R a set of strings;

if ($R \subseteq \Sigma$) return R ;

else:

choose a pivot string in $R \rightarrow P$

$$R_C = \{S \in R : S[i] < P[i]\}$$

$$R_ = \{S \in R : S[i] = P[i]\}$$

$$R_> = \{S \in R : S[i] > P[i]\}$$

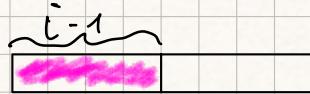
A = $MVK_Q(R_C, i)$;

B = $MVK_Q(R_ =, i+1)$;

C = $MVK_Q(R_>, i)$;

return A + B + C

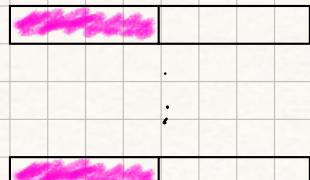
Fundament:



R

:

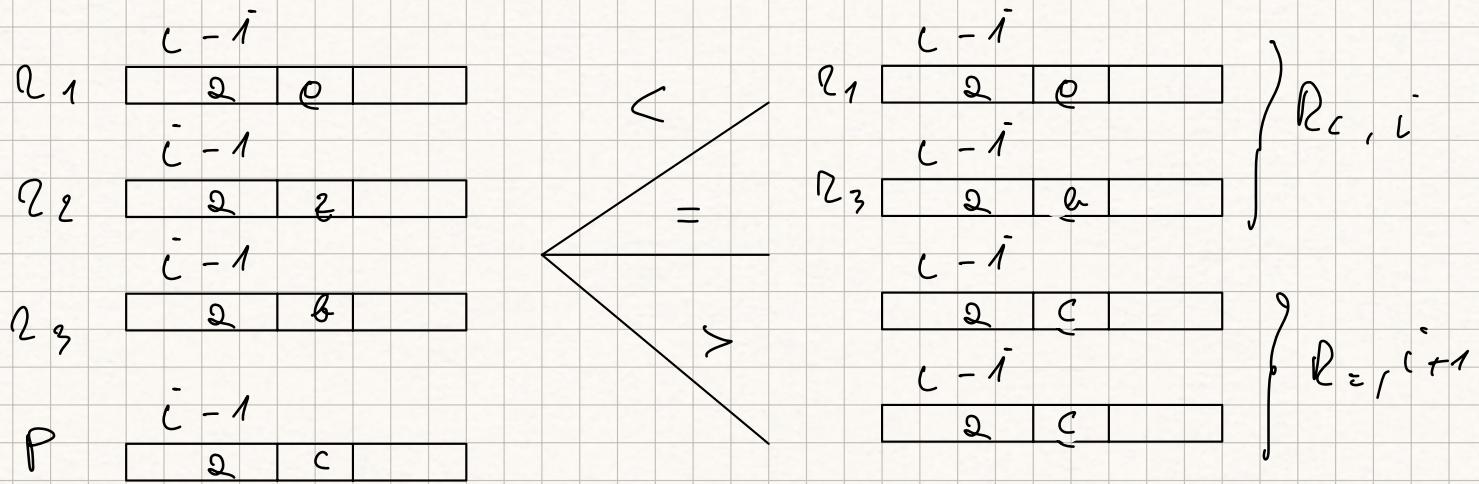
are equal



prefix free = no string in R

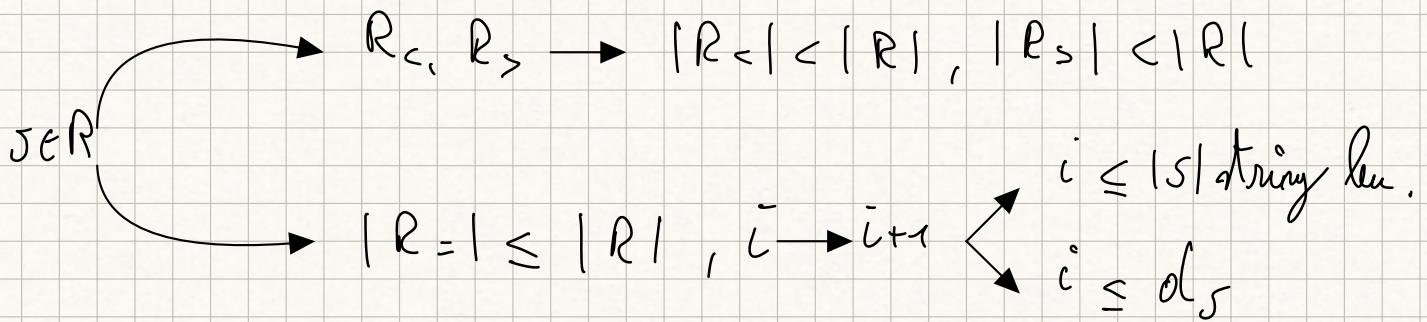
prefixes any other string in R

Assume that invariant holds



Assume that invariant holds (R, i)

Let consider a string $s \in R$



in the expected case, the pivot splits R in

a way that $|R_{<}| \cup |R_{>}| \approx \frac{|R|}{2}$

- $s \in R_{<} \cup R_{>} \text{, } O(\log n) \text{ times comparison}$

- cost of processing s is $O(\log n + o_s)$ time

$$\sum s O(\log n + o_s) = O(n \log n + o_l) \text{ lower bound}$$

$S = \{\text{cot}, \text{slai}, \text{cost}, \text{cor}, \text{st}\}$

$i = 1$ $P = \text{cost}$

$R = S$

$$R_C = \emptyset$$

$$P = AB1$$

$$R_{>} = \emptyset$$

$$R_{\leq} = \emptyset$$

$$R_{\geq} = \emptyset$$

$$P = AB1, AT$$

$$R_{\leq, 1} = \{AB1, AT\}$$

$$R_{\geq, 1} = \emptyset$$

$$P = AB1$$

$$R_{\leq} = \emptyset$$

$$R_{\geq} = \emptyset$$

$$R_{\leq, 2} = \{AB1\}$$

$$R_{\geq, 2} = \{AT\}$$

$$P = AB1, AT$$

$$R_{\leq, 3} = \{AB1\}$$

$$R_{\geq, 3} = \{AT\}$$

$$R_{\leq} = CAT$$

$$CATS$$

$$CAR$$

$$P = CAST$$

$$R_{\leq, 2} = \{CATS\}$$

$$R_{\leq, 3} = \{CAT\}$$

$$R_{\leq, 4} = \{CATS\}$$

$$R_{\leq, 5} = \{CAT\}$$

$$P = CAST$$

$$R_{\leq, 6} = \{CATS\}$$

$$R_{\leq, 7} = \{CAT\}$$

$$R_{\geq, 1} = \emptyset$$

$$R_{\geq, 2} = \emptyset$$

$$R_{\geq, 3} = \emptyset$$

$$R_{\geq, 4} = \emptyset$$

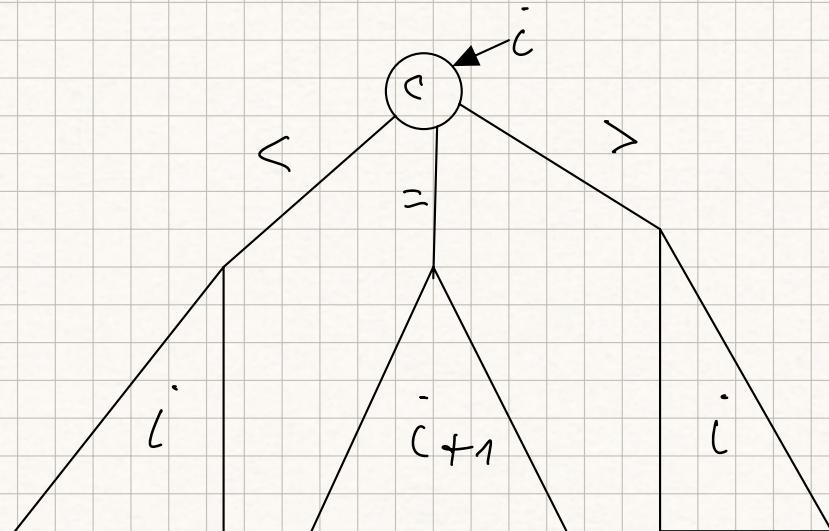
$$R_{\geq, 5} = \emptyset$$

$$R_{\geq, 6} = \emptyset$$

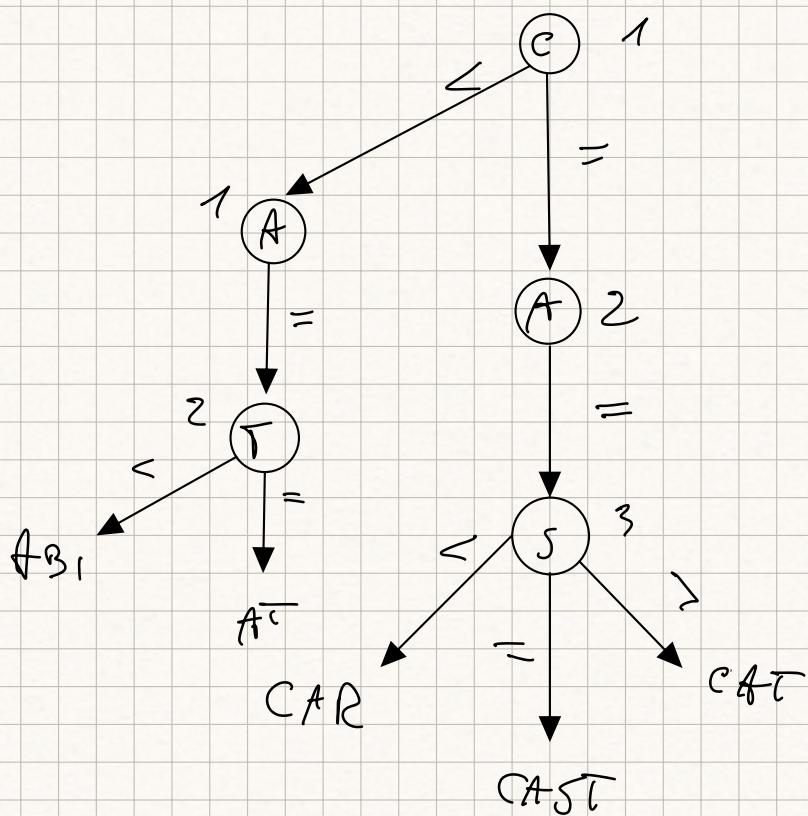
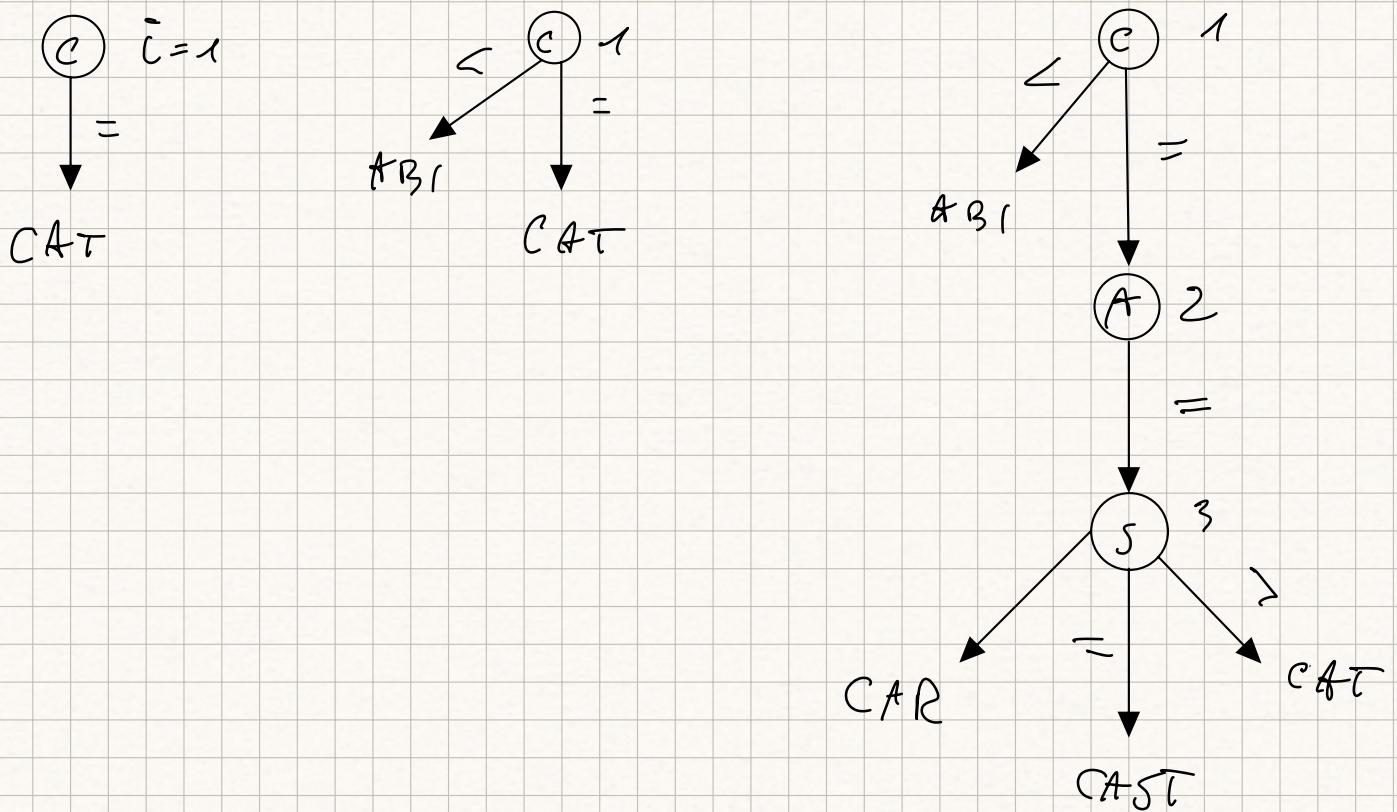
$$R_{\geq, 7} = \emptyset$$

$R_{>} = \emptyset$

Ternary Search Tree (TST)



$S = \{\text{cot}, \text{slai}, \text{cost}, \text{cor}, \text{st}\}$



If TST is balanced \rightarrow search $P : O(|P| + \log n)$

length of
searched string

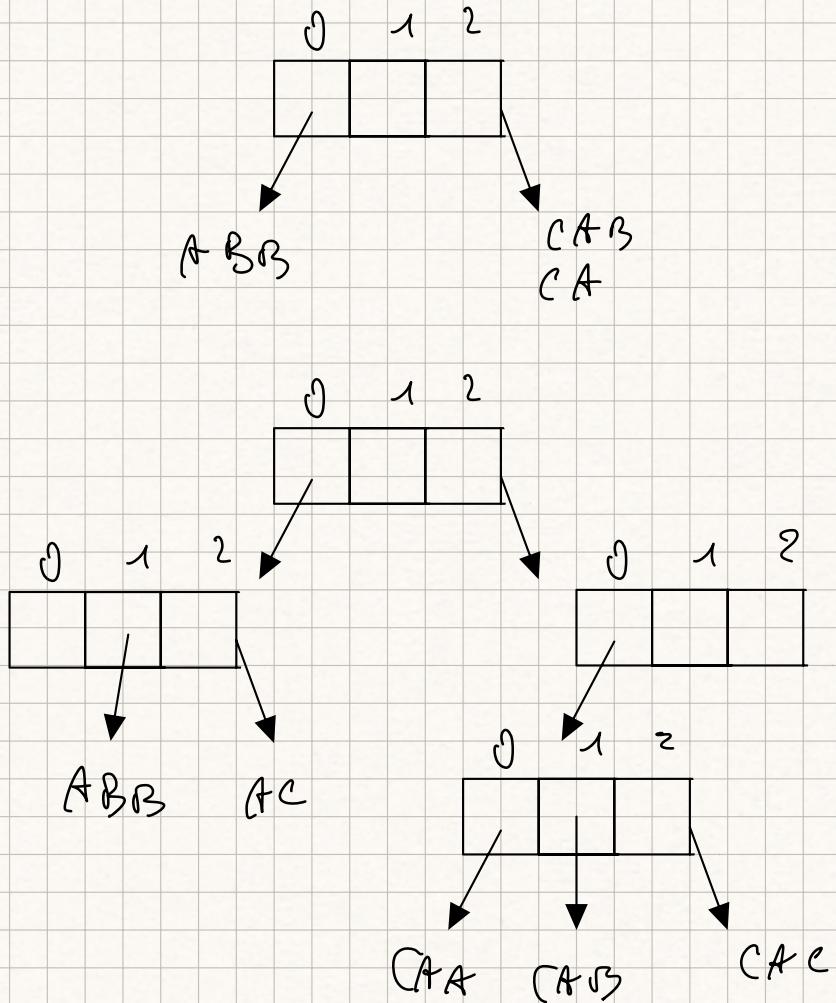
lower bound
of comparison tree

build a tree = multi-key insert = $O(N + m \log n)$

$$\Sigma = \{a, b, c\}$$

alphabet

trie, where hash-table is implemented by an array



$$\Sigma = \{A, B, C\}$$

- unweighted trie
- TST
- NKA [pivot is always the first string of RS]

Interpolation Search

| | | | | | | | | | | | | |
|-------|-------|---|-------|-------|---|-------|----|----------|----------|----|----------|----|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| $X =$ | 1 | 2 | 3 | 8 | 9 | 17 | 19 | 20 | 28 | 30 | 32 | 36 |
| | B_1 | | B_3 | B_6 | | B_8 | | B_{10} | B_{11} | | B_{12} | |

$n = \# \text{ Keys} \in N \approx x_1 < x_2 < \dots < x_n$ build once use many times

Query (y) $\rightarrow y \in X?$ \rightarrow binary search $O(\log n)$ time

Idee: Partition X in a "logical way" into buckets

B_1, B_2, \dots, B_m

biggest elem

smallest elem

dim of the universe

• Size of each bucket is $b = \frac{x_m - x_1 + 1}{m} = \frac{36 - 1 + 1}{12} = 3$

Portions of X

new array $I =$

| | | | | | |
|---|---|---|---|---|---|
| 1 | 3 | 0 | 0 | 4 | 5 |
|---|---|---|---|---|---|

$$|I| = 2n$$

$\underbrace{}_{B_1} \quad \underbrace{}_{B_2} \quad \underbrace{}_{B_3}$

build once use many times

of bucket = # of items

$$\text{Search}(y): I = \left\lfloor \frac{y - x_1}{b} \right\rfloor + 1 \quad O(1)$$

binary search in $B_I = X[I[2I-1], I[2I]]$

worst case $O(\log_2 \min\{n, b\})$

on expect. $O(\log_2 \text{count}) = O(1)$

The X obvious or random

$E[\max \# \text{Keys per bucket, when we have } n \text{ buckets}] = O(\log^2 n)$

$O(\log_2(\log^2 n))$ time w.l.o.g. = $O(\log(\log(n)))$ time

Time complexity is $O(\log_2 \Delta)$, where $\Delta = \frac{\max_{i=2..n} x_i - x_{i-1}}{\min_{i=2,3,\dots,n} x_i - x_{i-1}}$

for Interpolation search

measure of the distribution of the keys x_1, x_2, \dots, x_n

Disk:

$$\textcircled{1} \quad \max_{i=2, \dots, n} x_i - x_{i-1} \geq \sum_{i=2}^n x_i - x_{i-1}$$

$$= \frac{1 + x_n - x_1}{n} = b$$

$$\geq \frac{1 + \sum_{i=2}^n x_i - x_{i-1}}{n} =$$

medium

Telescopic sum

$$x_n - x_{n-1} + \\ x_{n-1} - \dots \\ \dots - x_1$$

$$\textcircled{2} \quad |B_S| \leq \frac{b}{\min_{i=2, \dots, n} x_i - x_{i-1}} \stackrel{\textcircled{1}}{\leq} \frac{\max_{i=2} x_i - x_{i-1}}{\min_{i=2} x_i - x_{i-1}} = \Delta$$

Set intersection $A \cap B$

A_m and B_m are sorted

Dict

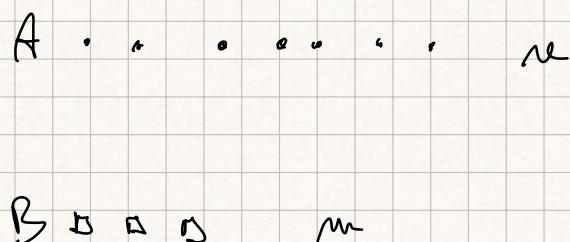
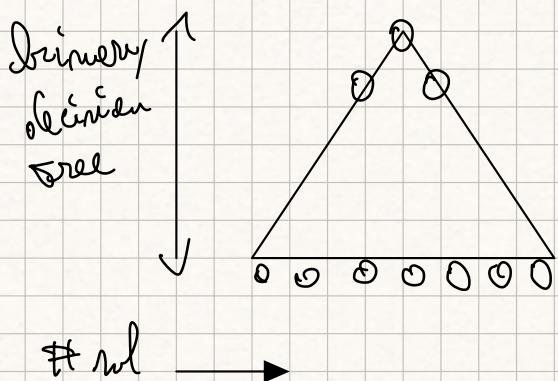
| | | |
|-------|---|-----------------|
| t_1 | → | 1, 3, 5, 11 |
| t_2 | → | 2, 4, 5, 10, 12 |

$m \leq n$

- Merge based: $O(n + m)$ Times optimal when $n \approx m$
- Binary search: $O(\min \log_2 n)$ Times i search the shorter in the longer list with B. S.

if: $\min \log n < n \rightarrow n < \frac{n}{\log n} \rightarrow$ Binary search optimal

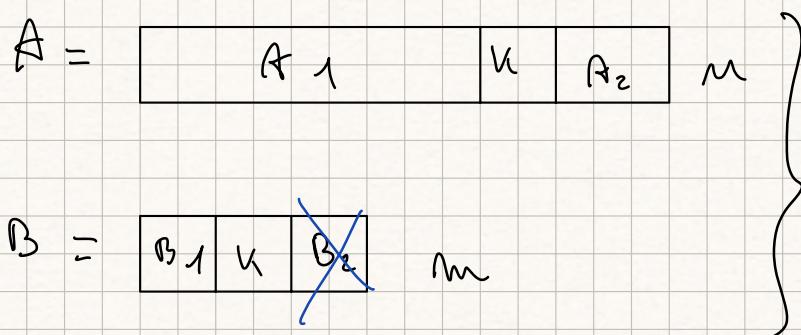
else: Merge search



$\binom{n}{m}$ halfs that
 $\binom{n}{m}$ B can be
included in A

$$\binom{n}{m} \# \text{comp} \geq \log_2 \# \text{al} \geq \log_2 \binom{n}{m} \approx \mathcal{O}\left(n \log_2 \frac{n}{m}\right)$$

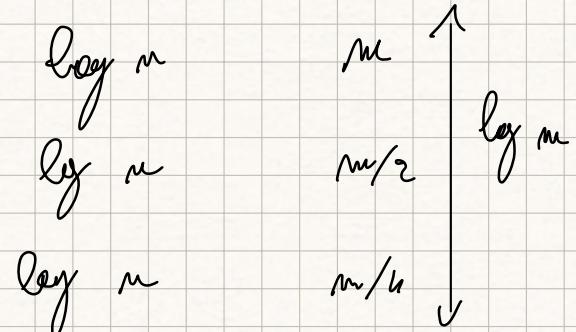
Mutual partitioning



$$A \cap B = \begin{cases} A_1 \cap B_1 \\ \vdots \\ A_k \cap B_k \end{cases}$$

① $(\log n)(\log m)$ time

unbalanced partition (left case)



$$\textcircled{2} \quad T(n, m) = O(\log(m) + 2T(n/2, m/2)) = O(m(1 + \log \frac{m}{n}))$$

balanced partition

$$m \leq n$$

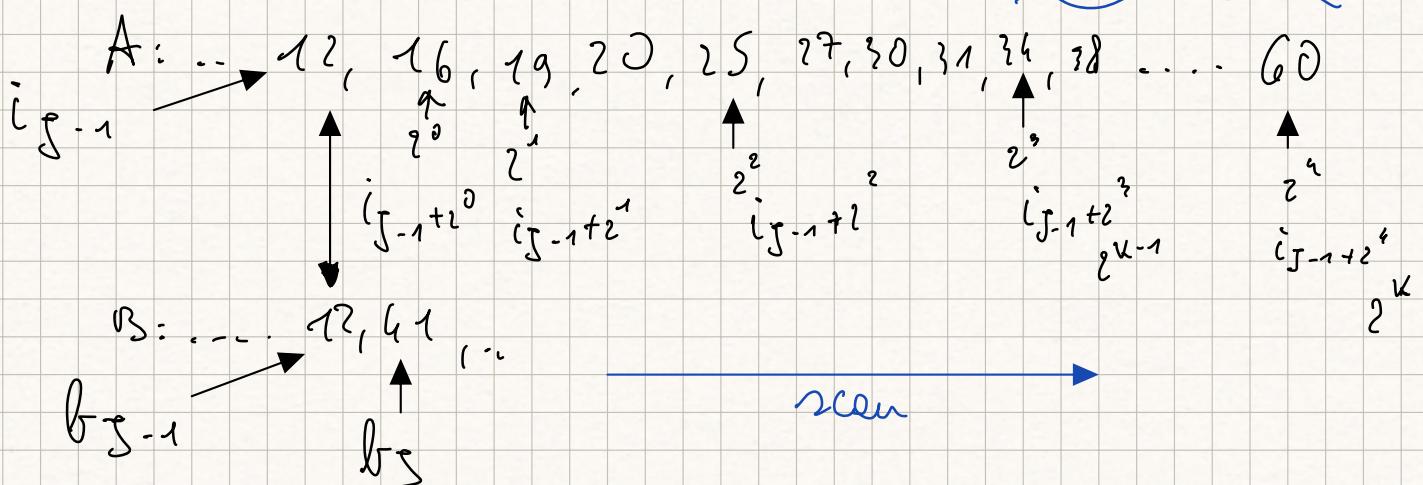
- $m \approx n \quad O(m) = O(n) \rightarrow$ as merge based
- $m \ll n \quad O(m \log n) \rightarrow$ as binary search based

LOWER BOUND

$$T(n, m) = \Omega(\log m) + 2T\left(\frac{n}{2}, \frac{m}{2}\right) = \Omega(m \log \frac{m}{n})$$

Doubling search, exponential search, ...

$$\Delta_j = \min(n, 2^{k-1}) \leq 2^{k-1}$$



$$\textcircled{1} \quad A[i_{j-1} + 2^{k-1}] < b_j \leq A[i_{j-1} + 2^k]$$

$$\textcircled{2} \quad 2^{k-1} < i_j - i_{j-1} \leq 2^k \quad \text{min}(i_{j-1} + 2^k, n) > i_j \geq i_{j-1} + 2^{k-1}$$

$$\textcircled{3} \quad \Delta_j \leq 2^{k-1} < i_j - i_{j-1} \leq 2^k$$

$$④ \log_2 2^{k-1} < \log_2 (i_j - i_{j-1})$$

$$k-1 < \log_2 (i_j - i_{j-1})$$

$$k = O(\log_2 (i_j - i_{j-1}))$$

$\leq k + k = O(k)$ steps we perform

$\xrightarrow{\text{B length}} O(\log_2 (i_j - i_{j-1}))$

$\xrightarrow{\text{Time}} = \sum_{j=1}^m O[\log_2 (i_j - i_{j-1})] = O\left(\sum_{j=1}^m \log_2 (i_j - i_{j-1})\right)$

$= O\left(\sum_{j=1}^m \log_2 (i_j - \bar{i}_{j-1})\right) \stackrel{\text{Jensen Inequality}}{=} \text{Telescopic sum}$

$$i_m - \cancel{i_0} \leq n$$

$= O(m \cdot \log_2 \frac{n}{m})$

Two-level storage approach

$$A = \begin{bmatrix} L \\ A_1 \\ L \\ A_2 \\ L \\ A_3 \\ L \\ A_4 \\ L \end{bmatrix} \quad m$$

$$A' = \begin{bmatrix} \square & \square & \square \end{bmatrix} \quad |A'| = \frac{m}{L}$$

$$B = \begin{bmatrix} \square & \square & \dots & \square \end{bmatrix} \quad m$$

$$1) A'[i+1] = A[1+iL] \quad \forall i=1,\dots,L$$

$$2) \text{ merge } (A', B) \Rightarrow B_i = \{b \in B : A'[i] \leq b < A'[i+1]\}$$

$$\therefore \forall i \quad A_i \cap B_i \quad i=1,2,\dots, \frac{m}{L}$$

$$T_1 = O(|A'| + |B|) = O\left(\frac{m}{L} + m\right) \text{ Time Merge expenses}$$

$$T_2 = \sum_{i=1}^{m/L} (|A_i| + |B_i|) = \underbrace{\sum_{i=1}^{m/L} |A_i|}_{\substack{\text{if } A_i \text{ picked} \\ \uparrow \\ m \cdot L}} + \underbrace{\sum_{i=1}^{m/L} |B_i|}_{m} \leq O(m + m \cdot L)$$

$$O(T_1 + T_2) = O\left(\frac{m}{L} + m \cdot L\right) \text{ Time}$$

Hashing \Rightarrow Dictionary problem

$$D = \{o_1, o_2, \dots, o_n\}$$

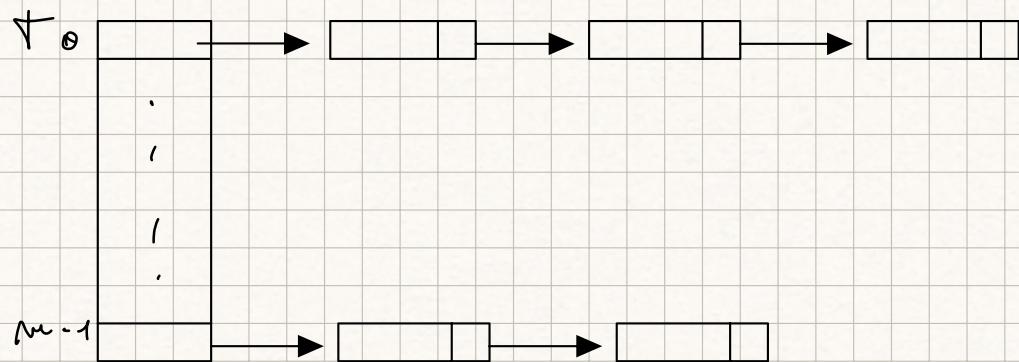
Key, satellite data

- Membership (v) / search (v)

- insert (o)

- delete (v)

Hashing with chaining



$$h : U \rightarrow [m]$$

collision resolution = chaining

Hashing useful iff $|U|$ is large

.) h is simple uniform 'hash function'

every key can go in every slot with uniform prob.

$$h(x) \begin{cases} 0 & 1/m \\ 1 & 1/m \\ \vdots \\ m-1 & 1/m \end{cases} \quad P(h(x)=i) \rightarrow 1/m$$

$$E[\text{dist } \tau[i]] = \sum_{k \in \Omega} k \cdot P(h(x)=i) = \sum_{k \in \Omega} \underbrace{\frac{1}{m}}_{1/m}$$

$$\text{load factor} = n/m = 2$$

$$n = \# \text{keys}$$

$$m = \# \text{hash entries}$$

• search compl. $\Theta(1) \Leftrightarrow \alpha = \Theta(1) \Leftrightarrow m \approx n \Rightarrow n = c \cdot m$

$$\text{search}(x) = \Theta(1 + \alpha)$$

$$m \geq n \rightarrow \alpha \leq 1$$

$$m < n \rightarrow \alpha > 1$$

• Global rebalancing technique
 n insertions

| | | | | | | |
|-----|----|----|----|---|---|----|
| m | 20 | 21 | 22 | . | . | 40 |
|-----|----|----|----|---|---|----|

when n is double

| | | | | | | |
|-----|----|---|---|---|---|----|
| m | 20 | . | . | - | . | 20 |
|-----|----|---|---|---|---|----|

of m is stop and it

| | | | | | | |
|-----|-----|--|--|--|--|---|
| 2 | 1 | | | | | 2 |
|-----|-----|--|--|--|--|---|

built $T' = 2 \cdot m$ and

ops.

i restart all the operat.

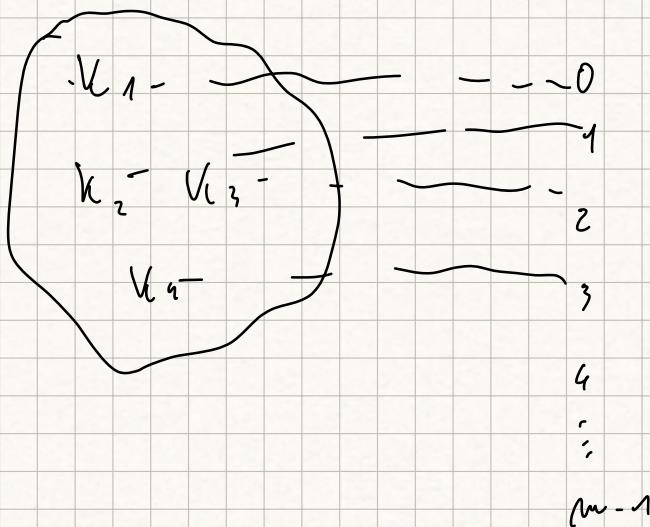
Amortized time complexity

uniform single hashing

$$h_i : U \rightarrow [m]$$

$\forall k \xrightarrow{h_i}$ any value in $[m]$

$\# h_i = \# \text{ all possible mappings from } U \text{ to } [m]$



| Vary | U_1 | U_2 | U_3 | ... | U_U | $\} \# \text{ mapping in total}$ |
|---------------|-------|-------|-------|-----|-------|----------------------------------|
| # ways to map | m | m | m | ... | m | |

is $m \cdot m \cdot \dots \cdot m = m^{|U|} \rightarrow \log_2 m^{|U|}$ bits

$$= |U| \cdot \log_2 m \text{ bits}$$

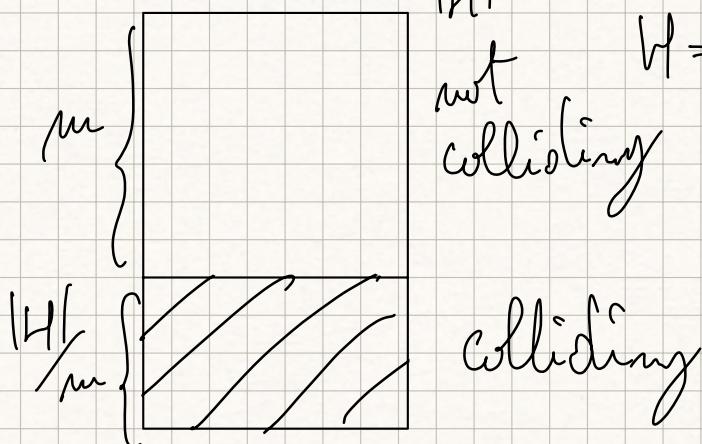
2^{64}

requires too much space

Universal hashing pairwise independence

The class of hash functions is called universal iff $\forall x, y \in U \quad \frac{|\{h \in H : h(x) = h(y)\}|}{m} \leq \frac{|H|}{m}$

$$P(h \text{ collides for } x, y) = \frac{|H|}{m} \cdot \frac{1}{m}$$



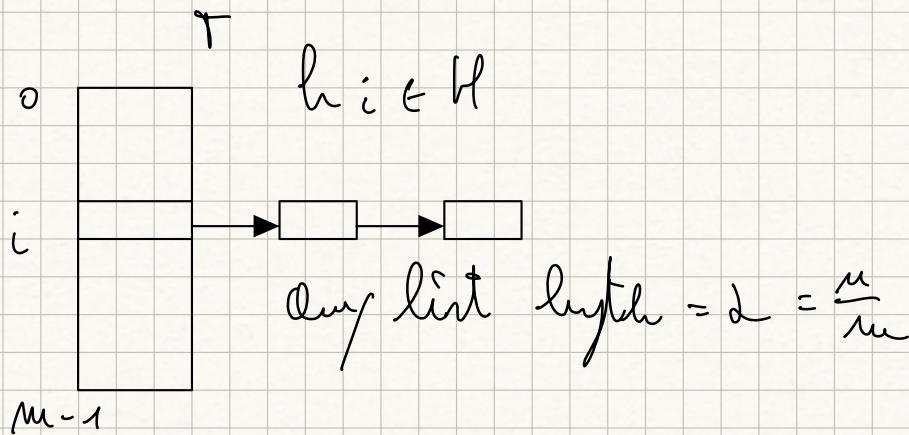
$|H|$

not colliding

$$H = \{h : U \rightarrow [m]\}$$

colliding

Class H exists \rightarrow Design hashing with chaining



Th: under the previous assumption, there exist one hash table with chaining of size m , in which search (k) over n objects takes $O(1 + \frac{n}{m})$ expected time.

Diss:

$$\mathbb{E} [\text{length } T[i]] = \sum_{k \in D} 1 \cdot P[h(k) = i] = 1 + \sum_{k \neq i} 1 \cdot P[h(k) = h(k')] =$$

Expected length of $T[i]$

$$k' \neq k; h(k') = i$$

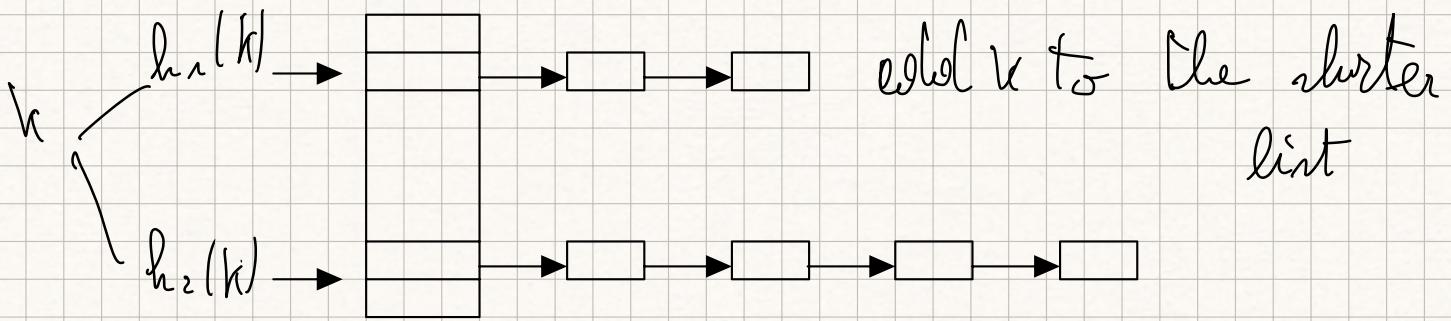
$$= 1 + \sum_{\substack{k \in D \\ k \neq k'}} \frac{1}{m} = 1 + \frac{m-1}{m} < 1 + \frac{m}{m} = 1+2$$

■

def of universal class

- expected len of a list in \mathcal{T} is $O(2)$

$$h_1, h_2 \in \mathcal{H}$$



$$e \text{ hashes} \rightarrow O\left(\frac{\log n}{\log \log n}\right) \quad \begin{matrix} \text{length of the longest list in } \mathcal{T} \\ \times 1 \text{ hash function} \end{matrix}$$

$$d \text{ hashes} \rightarrow O\left(\frac{\log \log n}{\log_2 d}\right) \quad \begin{matrix} \text{length of the longest list in } \mathcal{T} \\ \times d \text{ hash function} \end{matrix}$$

$$\mathcal{H} = \left\{ h_{\alpha}(v) = \sum_{i=0}^{n-1} \alpha_i \cdot v_i \bmod m, \text{ where } \alpha \in \{0, 1, \dots, m-1\} \right\}$$

$\alpha \in \mathbb{Z}$

Universal hash function

Prime

| | | | | |
|-----|-------|-------|---------|-----------|
| v | v_0 | v_1 | \dots | v_{n-1} |
|-----|-------|-------|---------|-----------|

$\log_2 n \text{ bits}$

$v_i \in \{0, 1, \dots, m-1\}$

α

| | | | |
|------------|------------|---------|----------------|
| α_0 | α_1 | \dots | α_{n-1} |
|------------|------------|---------|----------------|

$\log_2 m \text{ bits}$

$\log_2 |\mathcal{U}| \text{ bits}$

$\# \text{ of partitions for } v, \alpha$

$\alpha = \frac{\log_2 |\mathcal{U}|}{\log_2 m}$

$|\mathcal{H}| = |\mathcal{U}| - 1$

One function per value of α

Prove that \mathcal{H} is universal : means that, given $x \neq y$

$$\#\{h_{\alpha} : h_{\alpha}(x) = h_{\alpha}(y)\} \leq \frac{|\mathcal{H}|}{m} = \frac{|\mathcal{U}| - 1}{m}$$

Pick $x \neq y$

$x = \boxed{} \quad \boxed{} \quad \boxed{}$

$x_0 \quad \quad \quad x_{n-1}$

if

at least $x_0 \neq y_0$

$y = \boxed{} \quad \boxed{} \quad \boxed{}$

$y_0 \quad \quad \quad y_{n-1}$

if $h_{\alpha}(x) = h_{\alpha}(y)$ then $\sum_{i=0}^{n-1} \alpha_i x_i \bmod m = \sum_{i=0}^{n-1} \alpha_i y_i \bmod m$

$$\sum_{i=0}^{n-1} \alpha_i (x_i - y_i) \equiv 0 \pmod{m} \rightarrow \alpha_0(x_0 - y_0) + \sum_{i=1}^{n-1} \alpha_i (x_i - y_i) \equiv 0 \pmod{m}$$

$$Q_0(x_0 - y_0) \equiv - \sum_{i=1}^{n-1} Q_i(x_i - y_i) \pmod{m}$$

$$Q_0 \equiv -(x_0 - y_0)^{-1} \cdot \sum_{i=1}^{n-1} Q_i(x_i - y_i) \pmod{m}$$

$$\#\{\alpha : h_\alpha(x) = h_\alpha(y)\} \neq 0$$

exclude α_0

dim of $|U|$ for $\alpha = m^2$

$$\#\{\alpha \text{ that give a collision} : m \underbrace{-1}_{\neq 0} = \frac{m^2}{m} - 1 =$$

$$= \frac{|U|}{m} - 1 = \frac{|U| - m}{m} < \frac{|U| - 1}{m} \quad \blacksquare$$

$$h_\alpha, h_\beta(v_i) \equiv (\alpha x + b) \pmod{m} \quad \leftarrow \begin{matrix} \text{prime} \\ m > 1 \end{matrix} = ((\alpha x + b) \pmod{m}) \pmod{2^l}$$

$$h_\alpha(v_i) \equiv (\underbrace{\alpha \cdot k_i}_{2 \cdot \text{bits}} \pmod{2^l}) \text{ div } 2^{l-h}$$

$$m = 2^{l-20}$$

$$l \quad l-h$$

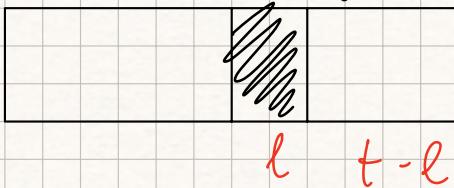
$$Q \cdot v = \boxed{}$$

$$h \quad h$$

scarto il primo blocco con $\text{mod } 2^h$
 scarto l'ultimo blocco con $\text{div } 2^{l-h}$

$$h(x) = (qx \bmod 2^t) \underbrace{\text{div } 2^{t-e}}_{\substack{\leftarrow \\ \in \mathbb{U}}} \rightarrow \text{integer in } [2^e]$$

$$h: \cup \rightarrow [\mu^n]$$



Cuckoo hashing:

Query $O(1)$ time

deletion $O(1)$ time

insertion $O(1)$ time amortized expected

Two choices: h_1, h_2

graph "logical"

Keys can move: relocated

| | | | | | | | | | | |
|---|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| + | A | C | | B | D | E | | F | | |
| | <u>2</u> | <u>1</u> |

let's insert G $h_1(G)$ $\overset{''}{0}$ $h_2(G)$ $\overset{''}{3}$

Th: If i, j positions in $[m]$, $\forall c > 1$, if $m \geq 2cm$

then $P(\exists \text{ undirected shortest path } i \dots j) \leq \frac{1}{m \cdot c^L}$

$\dim \text{ hash table} = \frac{m}{c} \leq \frac{1}{2c} c \frac{1}{2}$

$\Rightarrow \exp \text{ small in } L$

$$\mu = \# \text{ edges}$$

\downarrow

$$\# \text{ keys}$$

$$m = \# \text{ nodes}$$

\downarrow

$$\text{table size}$$

dim hash table $\frac{m}{c}$ keys

$$d = \frac{m}{m} \leq \frac{1}{2c} c \frac{1}{2}$$

$L = \text{length of the shortest path} \geq 1$

Proof: by induction on L

$$L=1: P(i \dots j) = 1 \text{ edge}$$

$$P\left(\exists k: h_1(k)=i, h_2(k)=j\right) \leq \sum_{k \in \Omega} P\left(h_1(k)=i, h_2(k)=j\right)$$

$$h_1(k)=j, h_2(k)=i\right)$$

$h_{1,2}$ = class of universal hash function

$$\sum_{k \in \Omega} 2 \cdot 1 \cdot \frac{1}{m} \leq 2 \sum_{k \in \Omega} \frac{1}{m^2} = \frac{2^m}{m^2} \leq \frac{m}{2c} \cdot \frac{2}{m^2} = \frac{1}{mc} = \frac{1}{m \cdot c^L}$$

$$L \geq 1: P(i \dots j) = P(\exists z: i \xrightarrow{\text{S.P.}} z \xrightarrow{L-1} j) \leq \sum_{\substack{\text{nodes } z \\ (|z|=m)}} P(i \xrightarrow{L-1} z \xrightarrow{1} j) =$$

$$= \sum_{z} \frac{1}{m \cdot c^{L-1}} \cdot \frac{1}{m \cdot c} = \sum_{z} \frac{1}{m^2 \cdot c^L} = \frac{m}{m^2 \cdot c^L} = \frac{1}{m \cdot c^L}$$

$$n = \# \text{ keys} \quad |T| = n$$

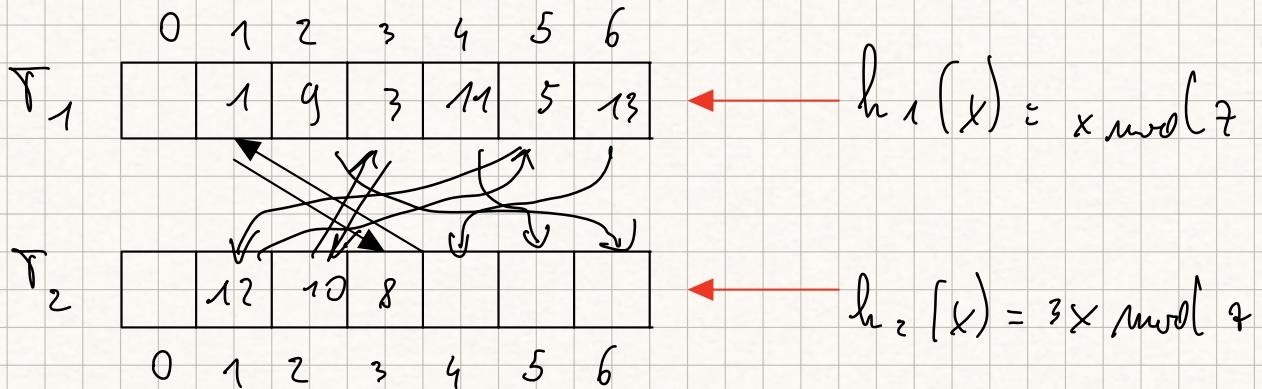
$$P(\text{head insertion}) = P(2 \text{ cycles}) \leq P(\geq 1 \text{ cycle}) = P(\exists i : \text{circle}_i)$$

$$P(\geq c, \exists i : \text{circle}_i) \leq \sum_{i=1}^m \sum_{l \geq 1}^{\infty} \frac{1}{m \cdot c^l} = \sum_{l \geq 1} \frac{m}{l \cdot c^l} = \sum_{l \geq 1} \left(\frac{1}{c}\right)^l = \frac{1}{c-1}$$

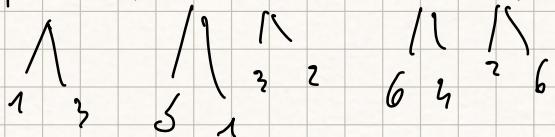
$$c = 3 \quad m \geq 2cm$$

Corollary: Given n keys and table of size m , if we take $m \geq 6n$ then $P(\text{head}) \leq \frac{1}{2}$

$$\hookrightarrow 2 \cdot \frac{n}{m} \leq \frac{1}{6} \approx 15\%$$



$$S = (1, 5, 8, 3, 17, 10, 11, 13, 9)$$



first: fill T_1 , then T_2 if $T_1[i]$ is full

Minimal Ordered Perfect Hashing

$$m = n$$

no collision

$$\forall v, v' \in D : h(v) \neq h(v')$$

$$\forall v' < v'' \in D : h(v') < h(v'')$$

rank(v')

(goal)

| $m=9$ | key | h | h_1 | h_2 |
|----------|-----|-----|-------|-------|
| elephant | 0 | 1 | 6 | |
| cat | 1 | 7 | 2 | |
| dog | 2 | 5 | 7 | |
| flop | 3 | 4 | 6 | |
| home | 4 | 1 | 10 | |
| house | 5 | 0 | 1 | |
| tom | 6 | 8 | 11 | |
| trip | 7 | 11 | 9 | |
| ton | 8 | 5 | 3 | |

prime number

① $h_1, h_2 : U \rightarrow [m']$ $m' \rightarrow m = m'$ universal
 universal
 $\begin{matrix} 13 \\ 9 \end{matrix}$ randomly chosen

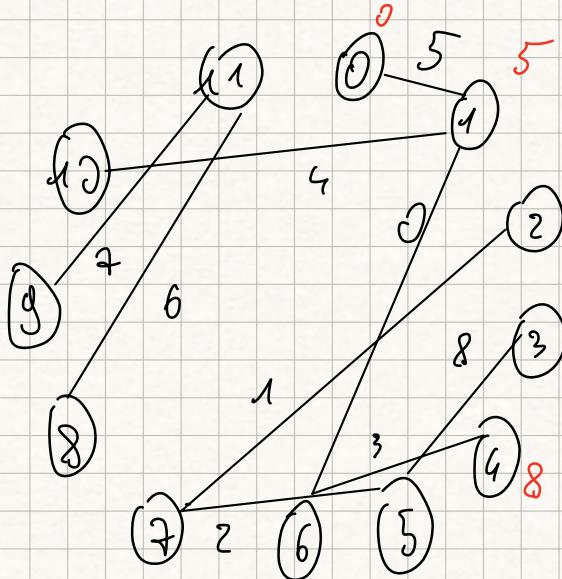
② $g(v) = [m'] \rightarrow [m]$ hash to be computed

We want to compute g so that

$$h(x) = [g(h_1(x)) + g(h_2(x))] \bmod n$$

We build G and
check if it is
acyclic → ok

No
restart from scratch



To define g it is:

$$h(\text{elecns}) = [g(h_1(\text{elecns})) + g(h_2(\text{elecns}))] \bmod 9$$

$$\textcircled{1} \quad 0 = g(1) + g(6) \bmod 9$$

$$\textcircled{2} \quad 1 = g(7) + g(2) \bmod 9$$

$$5 = g(0) + g(1) \bmod 9$$

We always start with the biggest connected component

$$\{0, 1, 6, 4, 10\} \setminus \{2, 7, 5, 3\} \setminus \{9, 11, 8\}$$

assign the first value to ^{given} or we prefer (typically 0)

| | | | | | | | | | | | | | |
|-------|---|---|---|---|---|---|---|---|---|---|----|----|----|
| $y =$ | 0 | 5 | 0 | 7 | 8 | 1 | 4 | 1 | 0 | 1 | 8 | 6 | 0 |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

$m' \cdot \log n + O(n)$ bits

$$\text{ls: } D = \{AA, BB, BD, CD\}$$

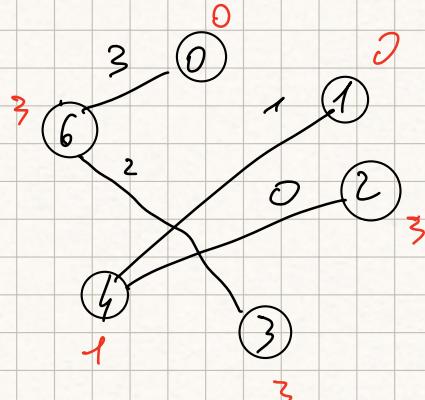
| | A | B | C | D |
|---|---|---|---|---|
| 2 | 1 | 2 | 3 | 4 |

$$h_1(x,y) = 3 \cdot h(x) + h(y) \bmod 7$$

$$h_2(x,y) = h(x) + 2h(y) \bmod 7$$

m' should be greater than # keys $2 > 4$

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



$\bmod 4$

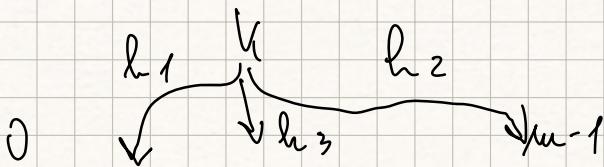
| | | | | | | | |
|-------|---|---|---|---|---|---|---|
| $g =$ | 0 | 0 | 3 | 3 | 1 | 0 | 3 |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 |

$$h(AD) = 0$$

$$\underbrace{g(h_1(AD))}_0 + \underbrace{g(h_2(AD))}_0 \bmod 4$$

Bloom Filter (Dictionary)

membership
insertion
~~deletion~~



① $B = \boxed{0 \dots 1 \dots 1 \dots \dots 1}$ bloom filter array $M > n$

② $h_1, h_2, \dots, h_r : () \rightarrow [n]$

$$\text{membership}(k) = B[h_1(k)] \wedge B[h_2(k)] \wedge B[h_3(k)] \wedge \dots \wedge B[h_r(k)]$$

$\begin{cases} 1 \text{ Yes } \approx \text{maybe} \\ 0 \text{ No for sure } \underline{\underline{100\%}} \end{cases}$

Keys are dropped after the insertion

$P(B[i] = 0 \text{ after inserting } n \text{ keys with 2 hash function})$

$$= \left(\frac{m-1}{m} \right)^{2 \cdot n} = \left(1 - \frac{1}{m} \right)^{2 \cdot n}$$

$$P(\text{error}) = P(y \notin D \text{ but } B[h_i(y)] = 1, \forall i=1, \dots, r) = \prod_{i=1}^r P(B[h_i(y)] = 1)$$

$$= \prod_{i=1}^r \left(1 - \left(\frac{m-1}{m} \right)^{2 \cdot n} \right) = \left[1 - \left(\frac{m-1}{m} \right)^{2 \cdot n} \right]^r = \left[1 - \left(1 - \frac{1}{m} \right)^{m \cdot \frac{2 \cdot n}{m}} \right]^r = e^{-1}$$

$$\approx \left[1 - e^{-\frac{2 \cdot n}{m}} \right]^r$$

Fixed m, m find \hat{r} s.t. $P(\text{error})$ is minimum

$$\Rightarrow \min_{\hat{r}} \left(1 - e^{-\frac{2 \cdot m}{\hat{r}}} \right)^2$$

$$\hat{r} = \frac{m}{m} \ln 2$$

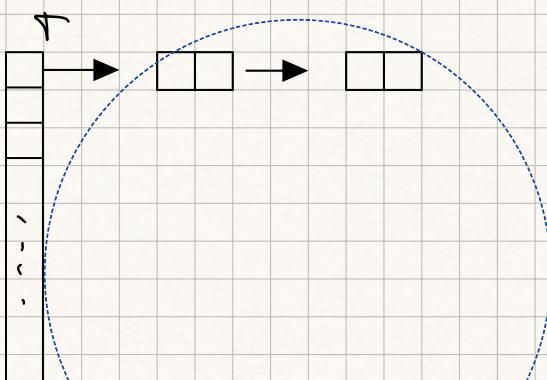
$$P(\text{error with } \hat{r}) = \left(\frac{1}{2} \right)^{\hat{r}} = \left(\frac{1}{2} \right)^{\frac{m}{m} \ln 2} = (0.6185)^{m/m}$$

① Fix $m \rightarrow \hat{r}$, guarantees minimum error for that m, m

② Fix $\hat{r}, \frac{m}{m} \rightarrow P(\text{error with } \hat{r}, \frac{m}{m})$
 ↑ speed ↑ memory

$$\frac{m}{m} \approx \# \text{ bits per key in BF}$$

$$\begin{aligned} m &= 10 \text{ M bits} \\ n &= 1 \text{ M keys} \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \frac{m}{n} = 10$$



$$\geq m \cdot \log n + n (\log |U| + \log n)$$

hashing with chaining: $m \cdot w + m(\log |V| + w)$

Let $m = c \cdot n$ where $c \geq 1$

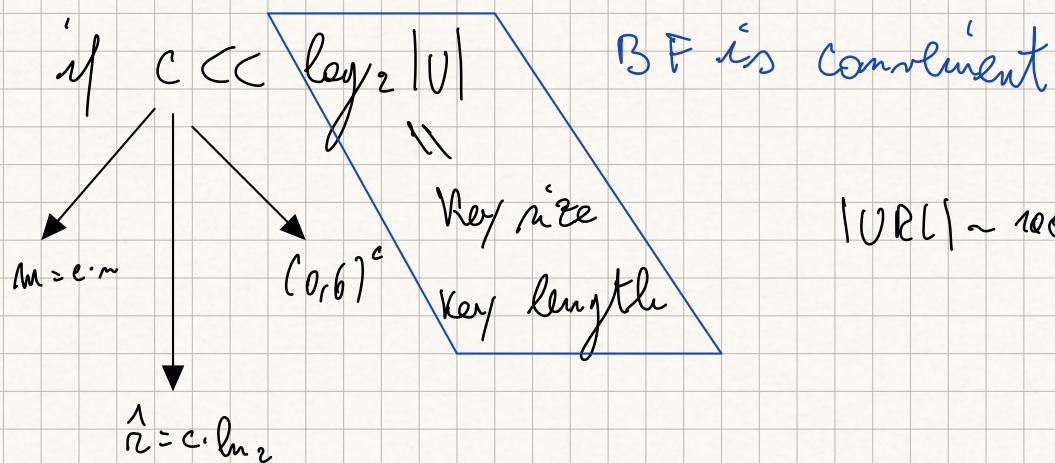
using c bits per key or storage

- Cuckoo hash $\rightarrow \frac{m}{n} \cdot \log |V| > \log |V|$ bits per key

- Hashing with chaining $\xrightarrow{m=n}$ ~~$\log n + (\log |V| + \log n)$~~ bits per key

$$= \log |V| + 2 \log n > \log |V| \text{ bits per key}$$

- BF takes c bits per key $\rightarrow P(\text{error}) \sim (0.6)^c$



$|URL| \sim 1000 \text{ chars} = 8000 \text{ bits}$

Define \hat{m} & B^* that minimizes $P(\text{error})$

$$m = 2^{20} \quad \hat{m} = 2^{40}$$

$$\hat{n} = \frac{\hat{m}}{m} \ln 2 = \frac{2^{40}}{2^{20}} \ln 2 = 2^{20} \ln 2$$

$$P(\text{error}) = (0.6185)^{2^{20}}$$

$\hat{n} = 10$ hash function

$$P(\text{error}) = \left(1 - e^{-\frac{2^m}{m}}\right)^2$$

$$\text{Find } m \text{ s.t. } P(\text{error}) \leq \frac{1}{1000} \quad \left(1 - e^{-\frac{2^m}{m}}\right)^2 \leq \frac{1}{1000}$$

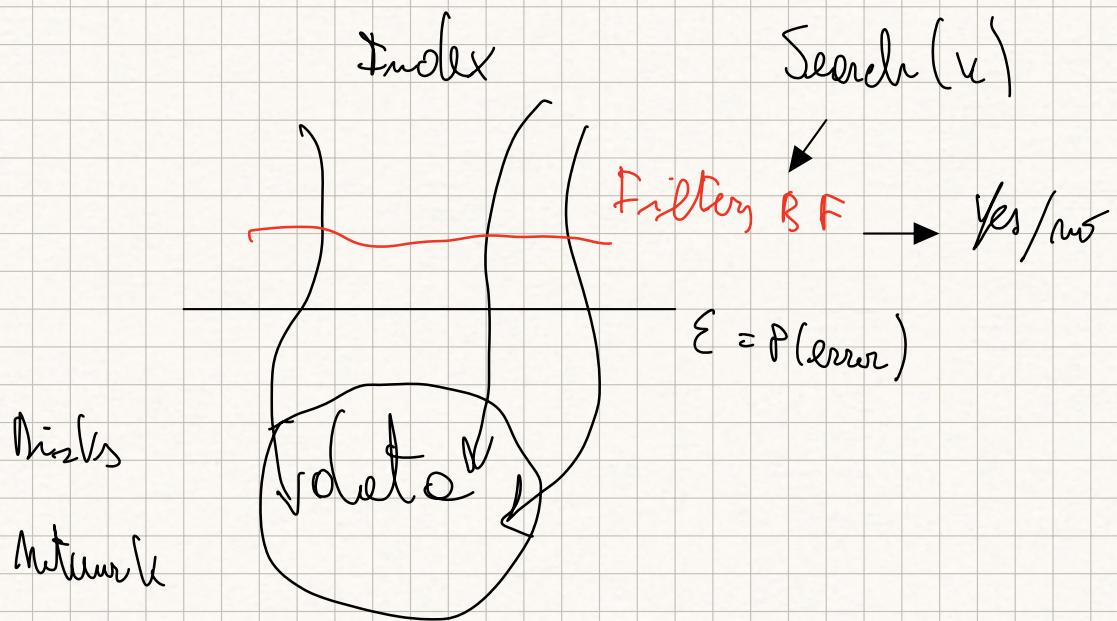
$$1 - e^{-\frac{2^m}{m}} \leq \sqrt[2]{\frac{1}{1000}}$$

$$e^{-\frac{2^m}{m}} \geq 1 - \sqrt[2]{\frac{1}{1000}}$$

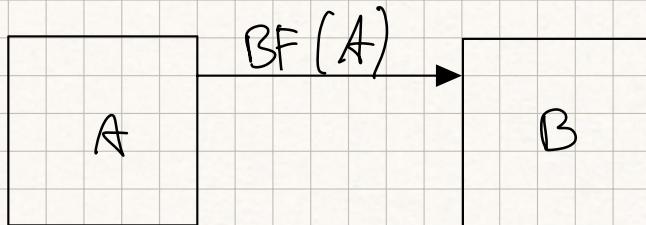
$$\frac{2^m}{m} \leq \ln \left(1 - \sqrt[2]{\frac{1}{1000}} \right)$$

$$m \geq \frac{\frac{2^m}{m}}{\ln \left(1 - \sqrt[2]{\frac{1}{1000}} \right)}$$

Rate bore (filters):



P2P (filters) : WVN (more volatile memory) $A \wedge B$



• $|A|$ log₂ $|V|$ bits

$BP(A)$ \leftarrow m_A bits
2 A hash function
 $m_A = |A|$

$\forall x \in B \rightarrow$ check if $x \in BF(A) \rightarrow YES(MAYBE)$

$$\text{MAP} : V \in B \rightarrow V \in A$$

$$V \notin A$$

$$(A \cap B) \cup (B - A)$$

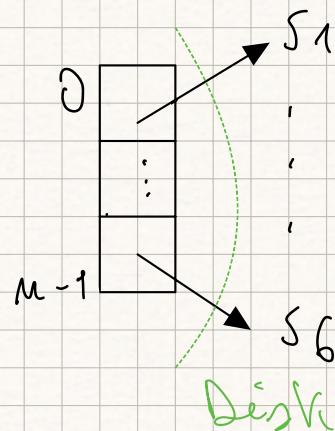
$$G[B-A] \leq P(\text{error}) \cdot |B| = (0.6185)^{\frac{m}{n^4}} \cdot |B|$$

Prefix (string) search: given a dictionary of n strings of total length N build a data structure so that searching $P[1, p]$ as a prefix of the dictionary strings is fast

$$D = \{S_1, S_2, S_3, S_4, S_5, S_6\}$$

$$S_1 = 000110, S_2 = 0001110, S_3 = 0010, S_4 = 0011, S_5 = 101, S_6 = 111$$

$P = 0001 \rightarrow$ 2 lexicographic searches $\rightarrow P\$$ ($\$ = \text{smaller than any symbol}$) $\rightarrow P\#$ ($\# = \text{bigger than any symbol}$)



exact binary search = $O(P \cdot \log n)$ time

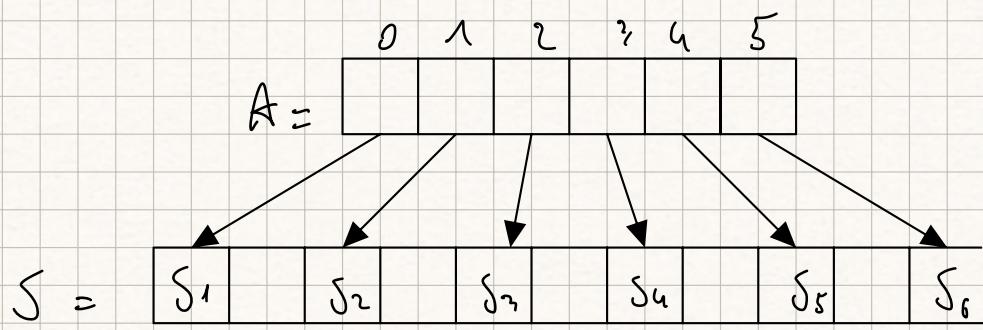
$$O\left(\frac{P}{B} \cdot \log n\right) \text{ I/Os}$$



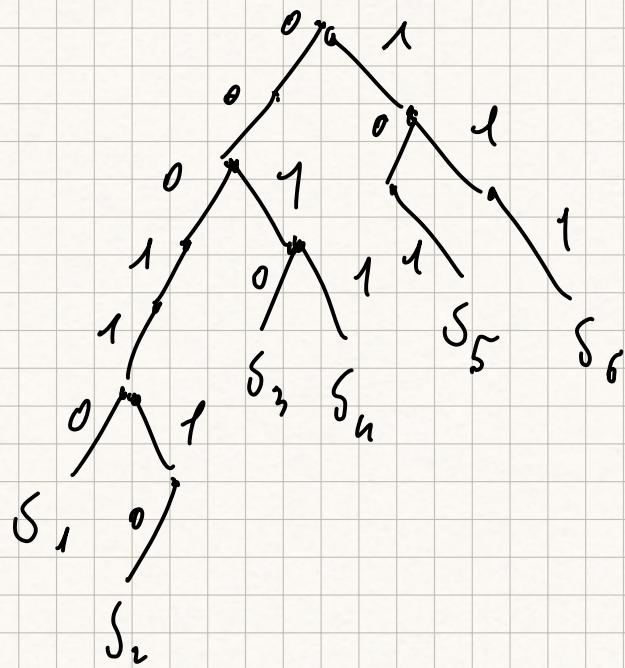
$$\text{Space} = N + (1 + w) m \text{ bytes}$$

constant size strings $\#$ strings

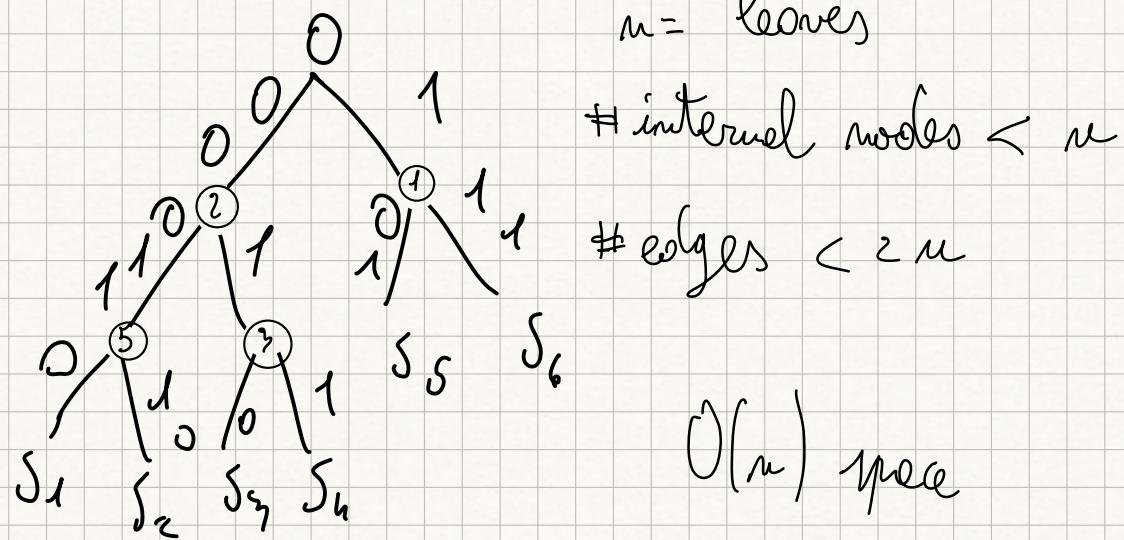
$\#$ partitions



Uncompressed trie



Compressed trie:

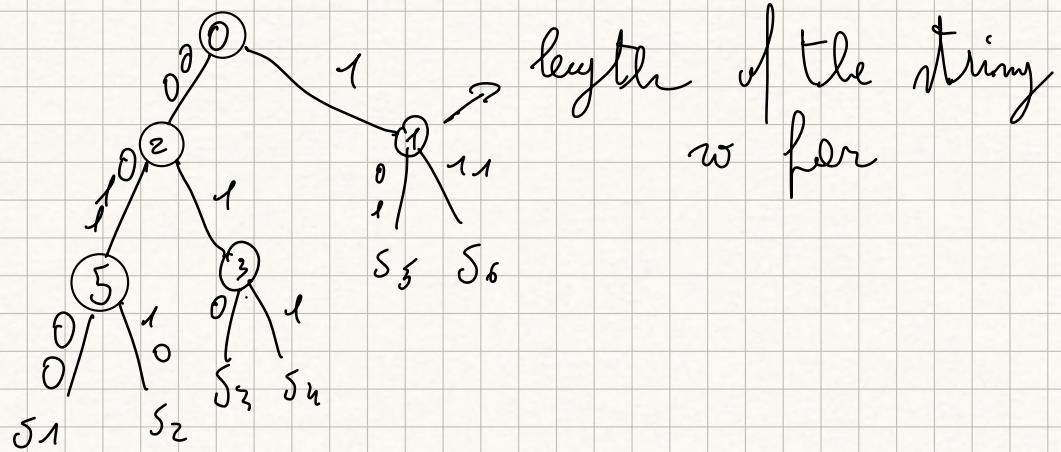


Compressed trie stored as: < string pointer, start, end >

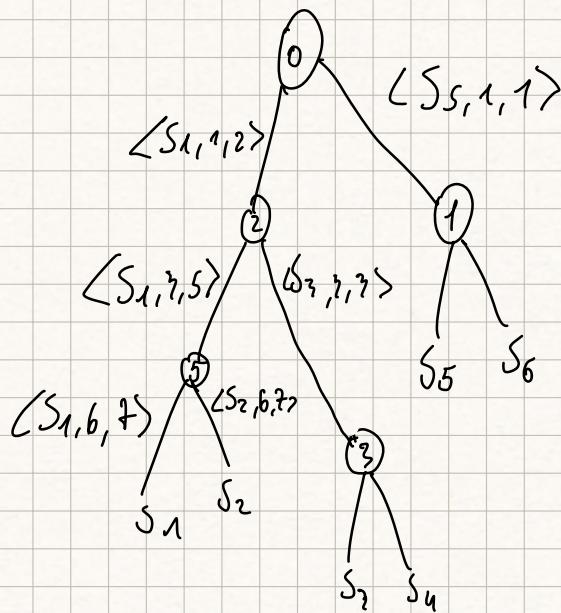
$(s_1, 1, 2) ; (s_2, 3, 5) \dots$

Patricia Tree:

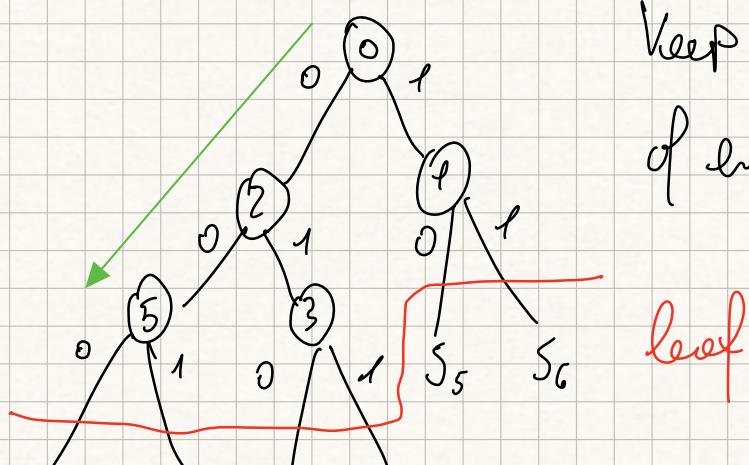
$$S = \{00011000, 0001110, 0010, 0011, 101, 111\}$$



In order to fit in internal memory we use pointers:



Patricia tree:



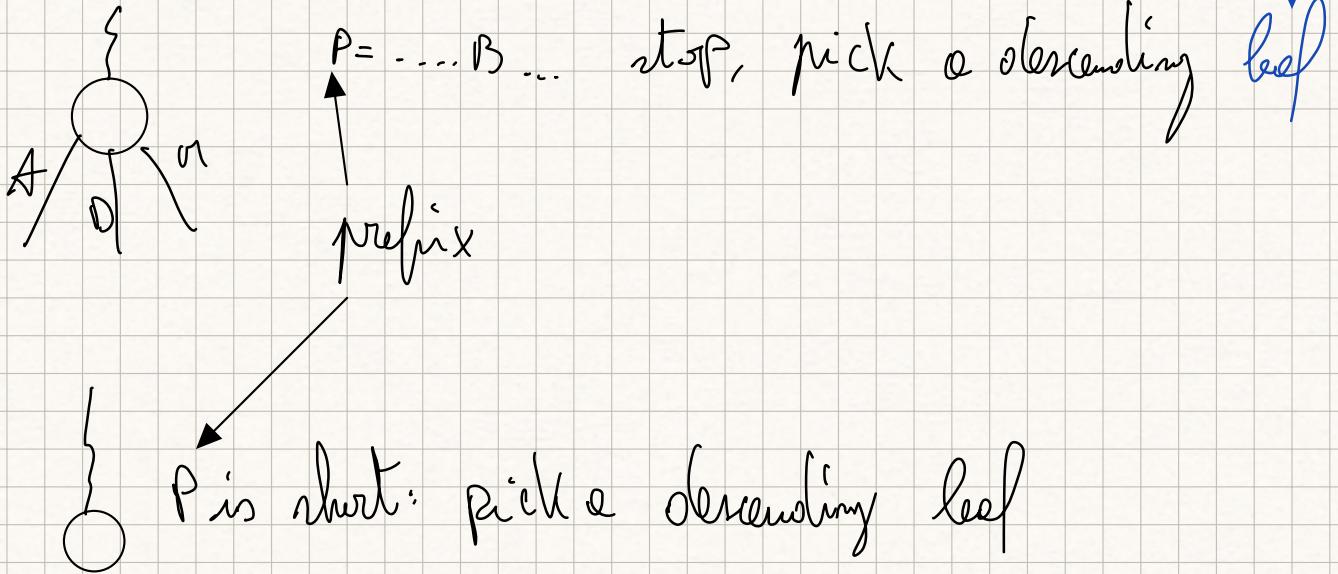
Keep only first edge
of every edge

leaf

$s_1 \quad s_2 \quad s_3 \quad s_4$

lexicographic search: $P = \underline{0} \underline{1} \underline{00} \underline{00} \underline{00}$

① Document search: match P's chars as much as possible



② Compute the longest common prefix between P and the picked leaf. Execute $O\left(\frac{P}{B}\right)$ I/Os
 Execute $O(P)$ time

③ upward traversal ↑

| | | | | |
|-----------|-----------|---------------|---------------|-------------|
| $s =$ | AA B A A | A B A B A A A | A B A B A B B | B B A B A |
| | A A B A B | A B A B A A B | A B A B B B B | B B A B A A |
| | A A B B A | A B A B A B A | A B B A A A A | B B A B A B |
| | B_1 | B_2 | B_3 | B_4 |
| $m = 1.5$ | | | | B B B A |
| | | | | B B B A B |
| | | | | B B B B B |

FC-storege per block

$\text{ABAAB} < 4, B \rangle, \langle 3, BA \rangle$

$\text{ABAABAAB} < 6, B \rangle, \langle 5, BA \rangle$

ABAABA

ABAABAAB

ABAABAAB

BBABAB

BBBA

DISK: B_1

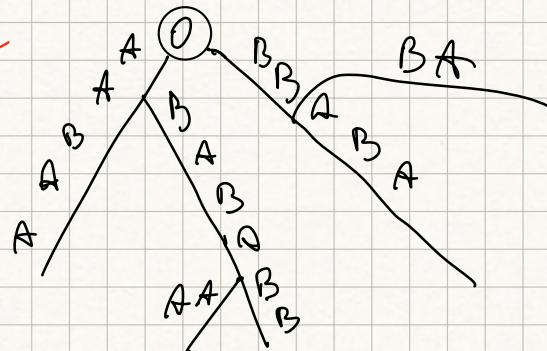
B_2

B_3

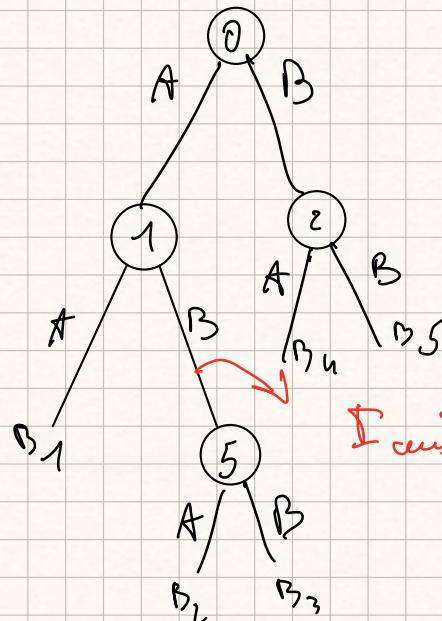
B_4

B_5

unsorted tree



Patricia Tree



I will reach between

B_3, B_4

01

P = AB BABB

Searching (Dictionary problem)

- Exact match : hash

- prefix match: unweighted/completed tries

Patricia tree

reduce substring to prefix
search

array of pointers

- substring match: suffix array/tree

Given a text $T[1, n]$ drawn from an alphabet Σ , preprocesses T and builds a d.s. so that, given on-the-fly a pattern $P[1, p]$ find all occurrences of P in T .

Count

list

retrieve

$$\cdot C = \{T_1, T_2, \dots, T_n\} \rightarrow T = T_1 \$ T_2 \$ T_3 \dots \$ T_n$$

Reduction from substring to prefix match

$$SUF(T) = \{T[i, n], \forall i: 1, 2, \dots, n\} \text{ set of all suffixes of } T$$

$$T = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ A & B & A & C & A & B & C \end{matrix}$$

$$SUF(T) = \left\{ \begin{matrix} A B A C A B C \\ B A C A B C \\ A C A B C \\ C A B C \\ A B C \end{matrix} \right\}$$

$$P = AB$$

↗ C
 C

If P occurs in T at position $i \Leftrightarrow P = T[i, i+|P|-1]$
 $\nwarrow P$ is a prefix of $T[i, m]$

Suffix array : array of pointers to text suffixes sorted
 alphabetically

| SA | |
|----|---------|
| 1 | ABACABC |
| 5 | ABC |
| 3 | ACABC |
| 2 | BACABC |
| 6 | BC |
| 4 | C |
| 9 | CABC |

$O(n)$ space

Data structure : SA and T
 $n \log n + n \log |\Sigma|$ bits

$$T = ABACABC$$

$P = AB$ # smaller than any character
 AB # longer than any character

$$\$ < \text{any } \Sigma < \#$$

$$O(n \log n) \text{ bits} = (n = \text{pointers} \cdot \log n = \text{any suffix len})$$

$$O\left(\frac{P}{n} \log n\right) I/O_s$$

- Count (P): prefix search from $P\$$ and $P\#$: $O(P \log n)$ time
- Retrieve (P): prefix search from $P\$$ and $P\#$: scan the subarray between positions of $P\$$ and $P\#$ and print

$$O(P \log n + \Theta cc) \text{ time}$$

Binary search p times

$$O\left(\frac{P}{B} \log n + \frac{\Theta cc}{B}\right) / Q$$

SA-build(τ, n):

```
for ( $i = 1; i \leq n; i++$ )
     $SA[i] = \tau + i - 1$ ; }  $O(n)$ 
```

qsort($SA, n, \text{sizeof}(\text{char}*)$, suffix_comp); } $O(n^2 \cdot n)$ avg $O(n(\log n))$

return SA ;

Suffix-comp($\text{char}** a, \text{char}** b$):

return str_cmp(*a, *b);

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

\uparrow
 avg # of compared chars

$T = ABBBCBBA$

Largest common prefix $[1, n-1]$

1 2 3 4 5 6 7 8

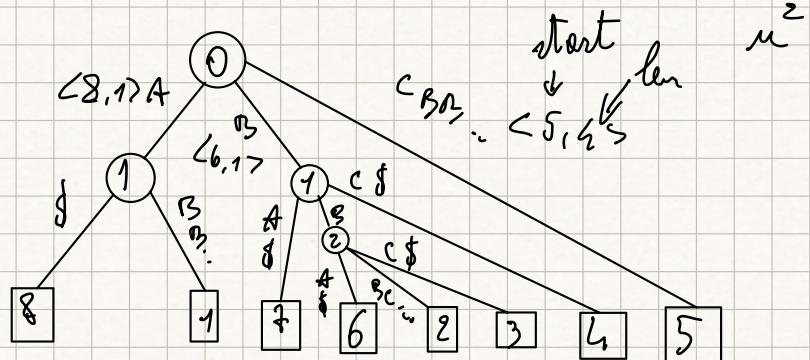
LCP SA

$\text{LCP}[i] = \text{largest common prefix}(\text{suffix } SA[i], \text{suffix } SA[i+1])$ for $i < n$

| | |
|---|-------------|
| 8 | A \$ |
| 1 | ABBBBCBBA\$ |
| 7 | BAB\$ |
| 6 | BBA\$ |
| 2 | BBBC\$ |
| 3 | BBC\$ |
| 4 | BC\$ |
| 5 | CBBBA\$ |

Suffix tree = compacted trie of all suffixes
of the string $T[1, m]$

SA + LCP \rightarrow ST



In a compacted trie the number of edges is proportional to the number of leaves

$$\begin{aligned} \# \text{ leaves} &= n \\ \# \text{ int. nodes} &< n \\ \# \text{ edges} &< n \end{aligned} \quad \left\{ O(n) \text{ space} \right.$$

With in order search we obtain the LCP.

With the leaves we obtain the SA.

- a) Check if it exists a substring of T that repeat ≥ 2 time
- We use LCP. We check ST

?) Check if it exists a substring of T that repeat ≥ 2 time

and has len = L

Check if an internal node of $S\Gamma$ has value $\geq L$

Check if LCP array has a value $\geq L$

?) Check if it exists a substring of T that repeat $\geq c$ time

and has len = L

In $S\Gamma$ we search for a node with value $\geq L$ and a number of child $\geq c$

In LCP array: $\exists LCP[i, i + c - 1] \geq L$

$c-1$

Diam:

$$T = \underbrace{\text{ } \underset{i}{\boxed{a}} \text{ } \underset{j}{\boxed{a}} \text{ } \underset{k}{\boxed{a}} \text{ } \underset{l}{\boxed{a}} \text{ } \dots}_{\geq L} \quad c = 3$$

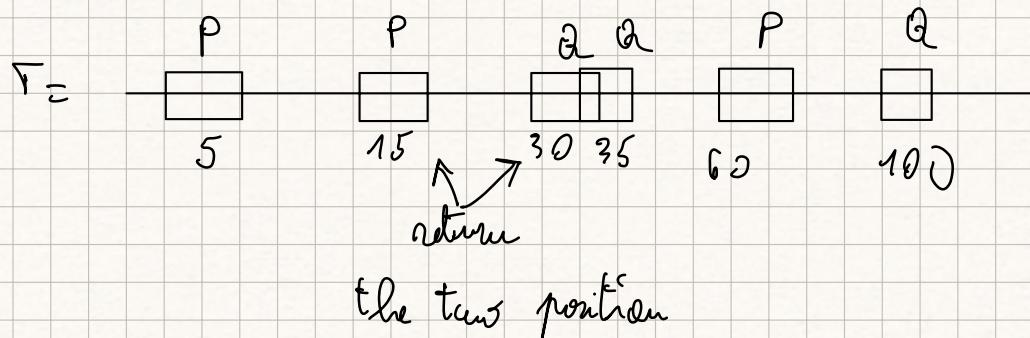
$\geq L$

$\geq L$

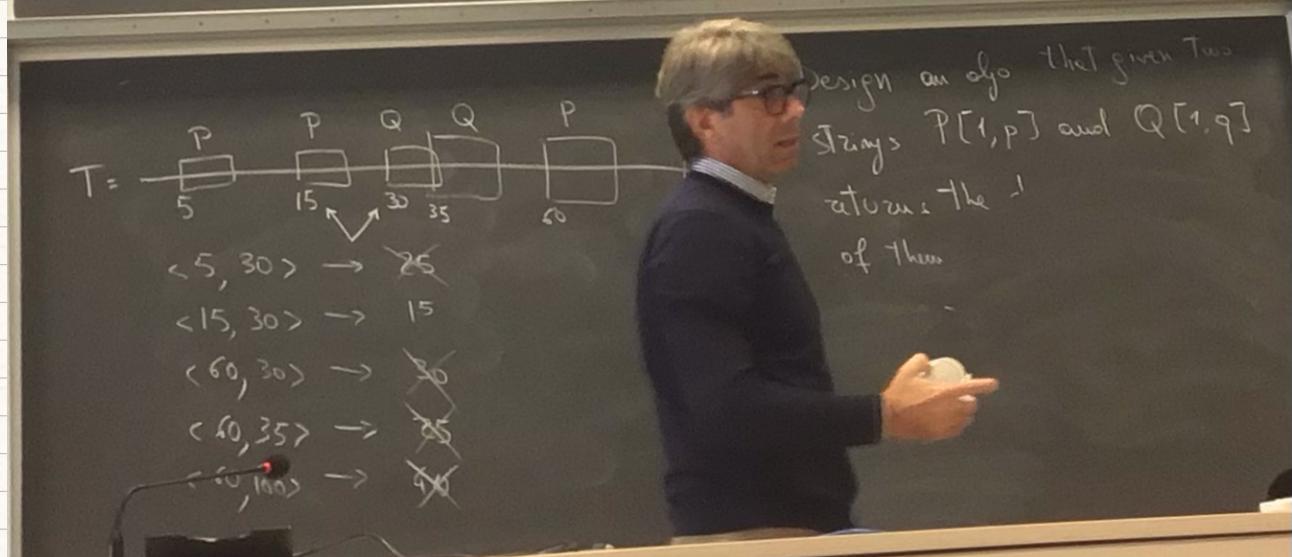
Design an alg. that given two strings $P[1, p]$ and $Q[1, q]$

return the closest occurrence of them in a text $T[1, n]$

Estimate the space and time complexity.



- 1) Search for $P \rightarrow L_p = \{5, 15, 60\}$ $O(p \log n)$ Time } $O(n)$ iff ST
- 2) Search for $Q \rightarrow L_q = \{30, 35, 100\}$ $O(q \log n)$ Time
- 3) $\text{Sort}(L_p) = \{5, 15, 60\}$ O_p
- 4) $\text{Sort}(L_q) = \{30, 35, 100\}$ O_q
- 5) Merge (L_p, L_q) Keeping the minimum absolute difference $O(O_p + O_q)$ Time } radix sort



$O(n)$ linear

Mutability

Source → emits symbols from an alphabet Σ

Probability P_f to emit Σ : $0 \leq P_f \leq 1$ $\sum P_f = 1$

information content of symbols

entropy of the source: $H_0 = \sum_f P_f \cdot \log_2 \frac{1}{P_f}$

no memory

measure the randomness of the source

$$\forall f \neq 0 \quad P_f \neq 0 \quad \text{iff} \quad \begin{cases} P_f = 0 \\ P_f = 1 \end{cases}$$

$$P_f = \frac{1}{|\Sigma|} \quad \leq \log_2 |\Sigma| \leq 1$$

$$\log_2 \frac{1}{P_f}$$

max entropy

$$\mathbb{E}[X] = \sum_f P_f \underbrace{\log_2 \frac{1}{P_f}}_{\text{log}_2 \frac{1}{P_f}}$$

$$X = \log_2 \frac{1}{P_f} \text{ with prob. } P_f$$

prefix-free codes

| | | |
|---------------------|---|---------------------------------|
| $0 \rightarrow 0$ | } | codewords (unique decodability) |
| $b \rightarrow 10$ | | |
| $c \rightarrow 111$ | | |

codewords are not prefix of each other

$$\begin{aligned} \Sigma &= \{0, b, c\} & H &= P(0) \cdot \log_2 \frac{1}{P(0)} + P(b) \cdot \log_2 \frac{1}{P(b)} + P(c) \cdot \log_2 \frac{1}{P(c)} \\ &\quad \frac{1}{3} \quad \frac{1}{2} \quad \frac{1}{4} & &= \frac{1}{3} \cdot \log_2 3 + \frac{1}{2} \cdot \log_2 2 + \frac{1}{4} \cdot \log_2 4 = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = 1.5 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{Average codeword length} &= \sum_f p(f) \cdot |\text{cw}(f)| = \\ &= P(0) \cdot 1 + P(b) \cdot 2 + P(c) \cdot 3 = \text{bits} \end{aligned}$$

$$= \frac{1}{4} + \frac{1}{2} \cdot 2 + \frac{1}{4} \cdot 3 = 1 + \frac{1}{4} + \frac{3}{4} = 2 \text{ bits/symbol}$$

The (Shannon): Any prefix-free code for some source

Σ takes on average # bits per symbol $\geq H$ (entropy)

$$\sum = N \quad X \in \Sigma \rightarrow \begin{cases} 2|x|+1 & \text{if } x < 0 \\ |x| & \text{if } x \geq 0 \end{cases} \in N$$

$S = S_1 S_2 S_3 \dots S_m$, $S_i \in N$, S_i may repeat

$S' = S_i \subset S_S$ increasing sequence

gap-excluding 1 2 11 3 ...

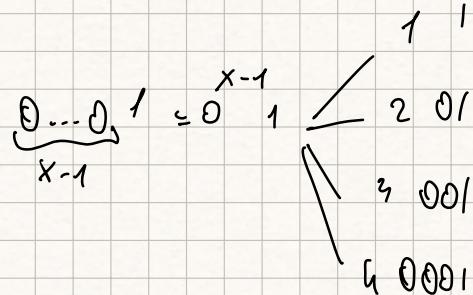
$$Q = \lceil \log_2(S^{*+1}) \rceil \text{ bits}$$

\downarrow

max S_i

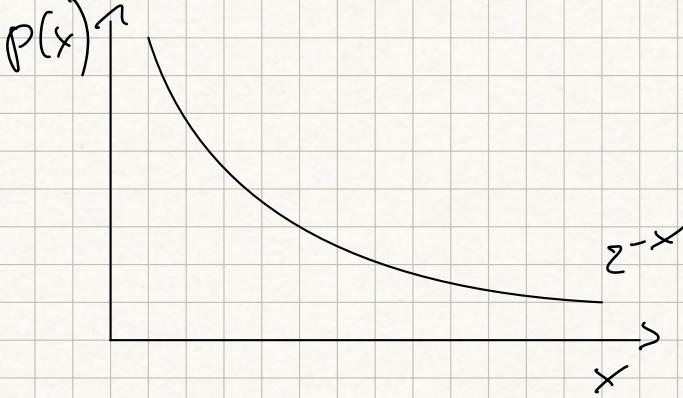
$S^* = 12 \quad S^{*+1} = 13$ fixed-length wordword code

Theory: $\forall x \geq 1$



If $|Cw(\sigma)| = \log_2 \frac{1}{P(\sigma)}$ code is optimal

$$|Cw(x)| = \log \frac{1}{P(x)} \Leftrightarrow x = \log_2 \frac{1}{P(x)} \Leftrightarrow 2^x = \frac{1}{P(x)} \Leftrightarrow P(x) = 2^{-x}$$



Elias:

$$\text{J-werk} : J(x) = \frac{0 \dots 0}{l-1} \frac{\text{bin}(x)}{l} \quad |J(x)| = 2^{l-1} = 2 \log_2 x$$

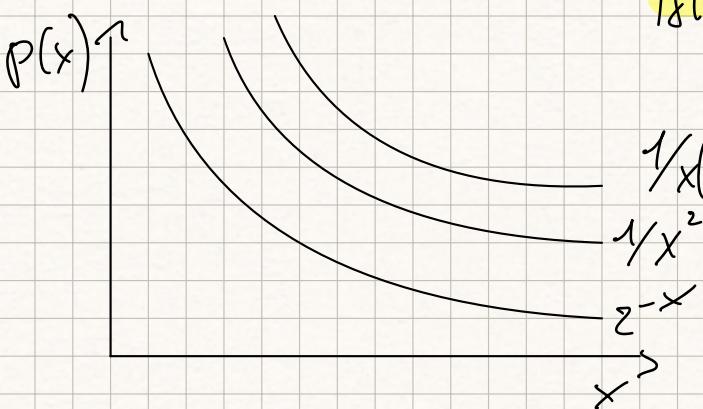
$$J(5) = 00101$$

$$|J(x)| = \log \frac{1}{P(x)} \Leftrightarrow 2 \log_2 x = \log_2 \frac{1}{P(x)} \Leftrightarrow P(x) = \frac{1}{x^2}$$

Pomer
low

$$f\text{-werk} : f(x) = \frac{J(l)}{2 \log l-1} \frac{\text{bin}(x)}{l} \quad f(5) = J(3)101 = 01101$$

$$|f(x)| = 2 \log_2 (\log x) - 1 + \log x \text{ bits}$$



$$2 \log(\log x) + \log x = \log_2 \frac{1}{P(x)}$$

$$\log(\log x)^2 + \log x = \log_2 \frac{1}{P(x)}$$

$$\log^2(x) \cdot x = \frac{1}{P(x)}$$

$P(x) \approx \frac{1}{\log^2(x) \cdot x}$

Variable byte:

$$x^0: 00 \underbrace{1\dots}_{\text{variable}}$$

bim
 1
 multiple of 7 bytes
 padding

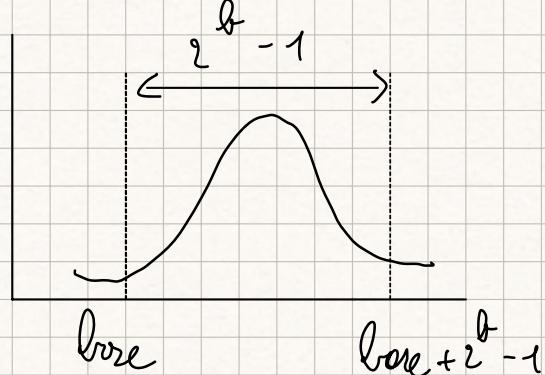
$$\left\lceil \frac{\text{bim}(x)}{7} \right\rceil \cdot 8 \text{ bytes}$$

tagging: add bit '1' in front of each block, '0' for the last one

$\text{bim}(x) = 100000011000101101010010$
 Padding Tagging

P for Delta: blocked scheme

Two parameters
base
 b



$x \in S : x \geq base$

partition S in two subsets

$$\begin{cases} \cdot x \in [base, base + 2^{b-1} - 1] = S_1 \\ \cdot x \geq base + 2^{b-1} = S_2 \end{cases}$$

$x \in S_1$, represent x in b bits

$$x - base \geq 0$$

$x \in S_2$: write $11\dots1$ and then encode x in a bytes

$$\begin{array}{c} \downarrow \\ 2^{b-1} \end{array}$$

exoptimal

$$S = \{0, 1, 6, 7, 3, 3, 9, 2\}$$

$$\{000, 001, 110, 111, 011, 011, 111, 010\}$$

$b_{\text{optimal}} = 0$

 $b = 3 \rightarrow [0, 6]$
 $\downarrow 3+1$

Take b such that $\geq 90\%$ of numbers in S are in S_1

Rice code (parameter k)

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

$$|R_k(x)| = \left\lfloor \frac{x-1}{2^k} \right\rfloor + 1 + k$$

$$R_4(20) = \begin{cases} q = \left\lfloor \frac{19}{16} \right\rfloor = 1 & O = 01 \\ r = 3 & \text{brim}_4(3) = 0011 \end{cases}$$

$$R_3(20) = \begin{cases} q = \left\lfloor \frac{19}{8} \right\rfloor = 2 & O = 001 \\ r = 3 & \text{brim}_3(3) = 011 \end{cases}$$

optimal distribution $P(x) = P(1-P)^{x-1}$

Blips - Few code

$S = S_1' S_2' \dots S_m'$ increasing (distinct)

$$\text{Miniblips} = S_m' + 1 \quad b = \lceil \log_2 m \rceil \quad l = \lceil \log_2 \frac{m}{n} \rceil \quad (\text{upper}) h = b + l$$

↑
(blipper) →
of numbers

| | |
|----|-------|
| 1 | 00001 |
| 4 | 00100 |
| 7 | 00111 |
| 18 | 10010 |
| 24 | 11000 |
| 26 | 11010 |
| 30 | 11110 |
| 31 | 11111 |

$$L = 0100111000101011 \quad |L| = n \cdot l = n \cdot \lceil \log_2 \frac{m}{n} \rceil$$

$$M = \underbrace{10}_{0}, \underbrace{110}_{1}, \underbrace{0}_{2}, \underbrace{0}_{3}, \underbrace{10}_{4}, \underbrace{0}_{5}, \underbrace{110}_{6}, \underbrace{11}_{7}$$

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16

$$|M| \leq 2m$$

negated every representation

for $(n+1)$ of every configuration

$$h = 3$$

$$b = 5 \text{ bits}$$

$$l = 2$$

$$n = 8$$

①

$$\#\text{1 in } M = m$$

$\#\text{0 in } M = 2^h$ $\#\text{0 is at most } n+1$ Represent M up to the
 in this example $2^h = m$ lost 1

② $\#\text{0 in } M \leq n$ (dropping the lost one)

$$\text{Total size} = n \left(\lceil \log_2 \frac{m}{n} \rceil + 2 \right) \text{ bits}$$

$$\text{Max } S = 18$$

$$n = 18 + 1 = 19$$

$$n = 8$$

$$b = \lceil \log_2 n \rceil = 5 \text{ bits} \quad l = \lceil \log_2 \frac{n}{m} \rceil = \lceil \log_2 \frac{19}{8} \rceil = 2$$

$$h = 3$$

| | | | | | |
|----|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 1 |
| 2 | 0 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 0 | 0 | 1 |
| 4 | 0 | 0 | 1 | 0 | 0 |
| 5 | 0 | 0 | 1 | 0 | 1 |
| 6 | 0 | 0 | 1 | 1 | 0 |
| 7 | 0 | 0 | 1 | 1 | 1 |
| 18 | 1 | 0 | 0 | 0 | 0 |

Access the i -th number:

① pick the i -th group of l bits in L

② $y = \underbrace{\text{select}_1(i)}_{\text{return position of}} - i$, represent y in h bits
 the i -th 1 in H

$$\text{Access}(3) = 001 \uparrow | 11$$

\uparrow
 3 -th compressed elem

$$y = \text{select}_1(3) - 3 = 4 - 3 = 1$$

$$\text{Access}(6) = 110 | 10$$

$y = \text{select}_1(6) - 6 = 12 - 6 = 6$

$\text{Rank}_1(i)$ } counting how many 1's are
 $\text{Rank}_0(i)$ } in the list until i -th pos.

$\text{Select}_1(i)$ } positions of
 $\text{Select}_0(i)$ } i -th 0/1

GEQ(i) \\\ greater or equal than i

$$p = \text{select}_0(H(i)) + 1;$$

$$j = p - H(i)$$

repeat :

$$y = \text{Access}(s);$$

$$j++;$$

$$\text{while } (y \geq i) \& (s \leq n)$$

$$es: 0, 3, 5, 8, 9, 10, 11$$

$$\text{Access}(s) = 9$$

$$GEO(6) = 8$$

$$l = \lceil \log_2 12 \rceil = 4$$

$$l = \lceil \log_2 \frac{12}{4} \rceil = 1$$

$$h = 3$$

$$\text{Access}(s) \Rightarrow \text{select}_0(s) - 5 = 9 - 5 = 4 \Rightarrow 100$$

$$GEO(6) \Rightarrow 6 = 0110 \quad p = \text{select}_0(3) + 1 = 2$$

$$j = 2 - 3 = 4$$

$$\text{Access}(4) \Rightarrow \text{select}_0(4) - 4 = 4 \quad 1000$$

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 0 |
| 1 | 0 | 1 | 1 |

$$L = 0110101$$

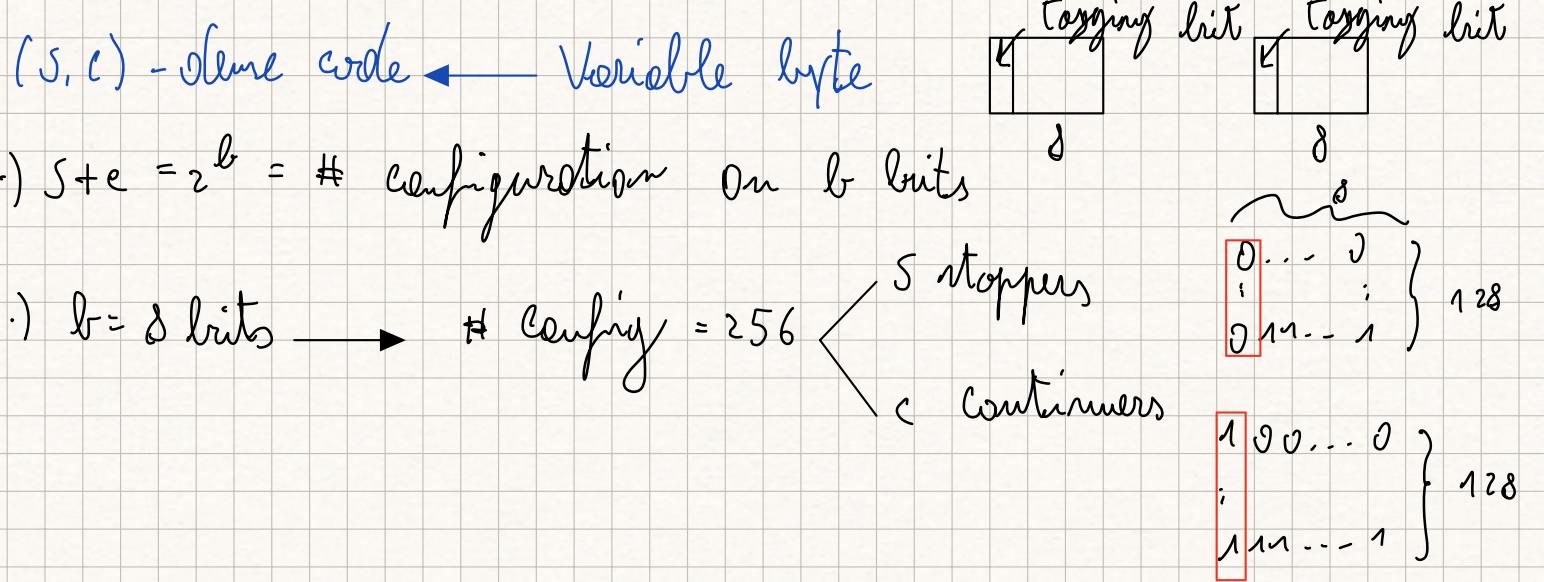
$$H = 1010100110110$$

Configurazione del primo
ottenere l'individuamento del
bucket della configurazione $H(\text{bin}(i))$

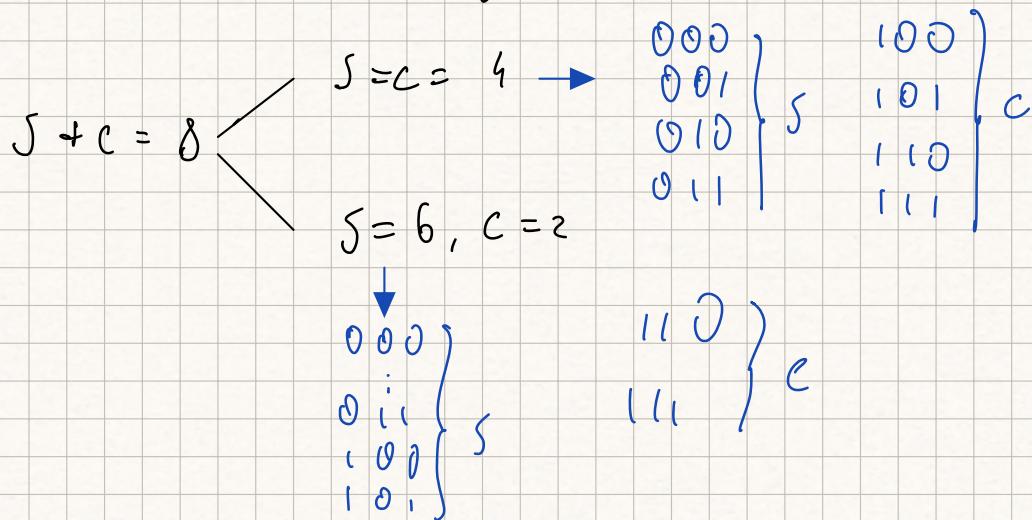
return the position of $H(\text{bin}(i))$ zeros
in bl of Elias-Fano

ottengo la posizione del minor elements
del bucket corrispondente alla $H(\text{bin}(i))$

del bucket corrispondente alla $H(\text{bin}(i))$

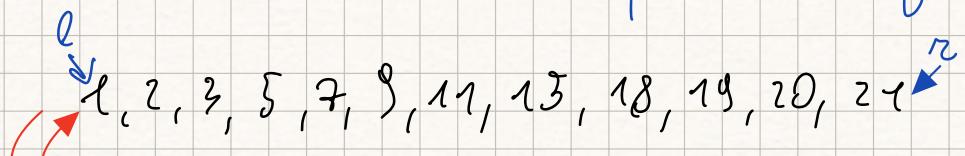


if $b = 3 \text{ bits} \Rightarrow 2^3 = 8 \text{ config}$



| Number | Code word (4,4) | Code word (6,2) | Code word (1,7) |
|--------|-----------------|-----------------|-----------------|
| 0 | 000 | 000 | 000 |
| 1 | 001 | 001 | 001 000 |
| 2 | 010 | 010 | 010 000 |
| 3 | 011 | 011 | 011 000 |
| 4 | 100 000 | 100 | 100 000 |
| 5 | 100 001 | 101 | 101 000 |
| 6 | 100 010 | 110 000 | 110 000 |
| 7 | 100 011 | 110 001 | 111 000 |
| 8 | 101 000 | 110 010 | C C S |
| | | | 001 001 000 |

Interpolation look (rep increasing)



$1, 1, 2, 2, 2, 4, 3, 1, 1, 1$ 1^{st} gap - encoding

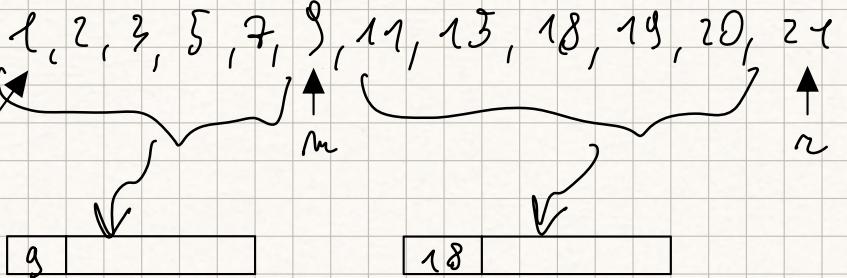
parameters : $l, r, \text{low}, \text{high}$

| | | | |
|-----|-----|--------------|---------------|
| l | r | low | high |
| $"$ | $"$ | $"$ | $"$ |
| 1 | 12 | 1 | 21 |

position

$$m = \left\lfloor \frac{l+r}{2} \right\rfloor = 6 \rightarrow S_6 = 9$$

using the minimum # bits l



left recursive cell : $\langle 1, 5, 1, 8 \rangle$

right recursive cell : $\langle 7, 12, 10, 21 \rangle$

| | | | |
|-------|-----|--------------|---------------|
| $m+1$ | 2 | low | high |
| l | | | |

$$Q = \text{low} + m - l \leq j \leq \text{high} - r + m = b$$

$$Q = 1 + 6 - 1 = 6 \quad b = 21 - 12 + 6 = 15$$

$$m = 6, S_m = 9 \rightarrow \text{low} + m - l = 6 = Q$$

| | | |
|-----|-----|-----|
| 1 | 6 | 1 |
|-----|-----|-----|

$$\rightarrow \text{high} - r + m = 15 = b$$

decode $S_m - \varrho = 3$ using $\lceil \log_2 \frac{15-6+1}{\varrho-\varrho+1} \rceil = \lceil \log_2 10 \rceil = 4$ bits

$\text{left } < l, r, \text{low}, \text{high} \Rightarrow \text{left } < l, m-1, \text{low}, S[m]-1 >$

$$\begin{array}{cccc} 1 & 5 & 1 & 8 \end{array}$$

$$m = 3, S_m = 3 \quad \left\lceil \frac{l+2}{2} \right\rceil = m$$

$$\varrho = 1 + 3 - 1 = 3 \quad \text{low} + m - l$$

$$b = 8 - 5 + 3 = 6 \quad \text{high} - r + m$$

decode $(S_m - \varrho) = 3 - 3 = 0$

in $\lceil \log_2 (6-3+1) \rceil = 2$ bits

$(00)_2$

$\text{right } < l, r, \text{low}, \text{high} \Rightarrow \text{right } < m+1, r, S[m]+1, \text{high} >$

$$\begin{array}{cccc} 7 & 12 & 10 & 21 \end{array}$$

$$m = 9, S_m = 18$$

$$\varrho = 10 + 1 - 7 = 12 \quad \text{low} + m - l$$

$$b = 21 - 12 + 9 = 18 \quad \text{high} - r + m$$

decode $S_m - \varrho = 18 - 12 = 6$

in $\lceil \log_2 (18-12+1) \rceil = 3$ bits

$(110)_3$

Compressed/computed sl. 3.

$\text{Rank}_b(i) = \text{count in } B[1, i] \text{ the number of } b \text{ bits}$

$\text{Select}_b(i) = \text{position of the } i^{\text{th}} \text{ bit} = b$

$\text{Rank}_1(i) = i - \text{Rank}_0(i)$

$b = [0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0]$

$$\begin{array}{c} \uparrow \\ \text{Select}_1(2) = 3 \end{array} \quad \begin{array}{c} \uparrow \\ \text{Rank}_1(5) = 2 \end{array}$$

$$\approx \frac{m}{\log_2 m}$$

th 1: I data structure that occupies $O(m)$ bits in addition

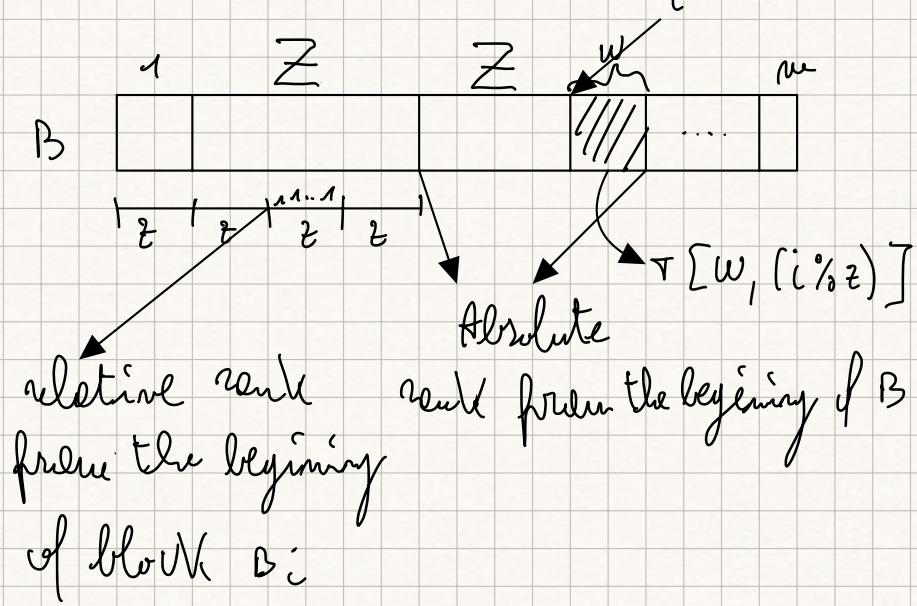
to $B[1, m]$ and report Rank_1 in $O(1)$ time

/Select₁ /Select₀
i

$$\# \text{Big blocks} = \frac{m}{z}$$

$$\# \text{Small blocks} = \frac{m}{z}$$

$$\text{Time} = O(1) + O(z)$$



relative rank rank from the beginning of B
from the beginning
of block B :

$$\text{Space in bits} = O\left(\frac{m}{z} \cdot \log m + \frac{m}{z} \cdot \log z + z \cdot 2^z \log z\right)$$

space of absolute rank
 of big blocks space of relative rank
 of small blocks space of lookup
 Table T

Four-routine Trick
 (configuration of small
 block)

$$\text{Rank}_1(w, z) = \#\text{1 in } w[1, z]$$

z
 configuration
 11...1

| T | 1 | 2 | . | . | . | z | position in z |
|-------|---|---|---|---|---|---|------------------|
| 00..0 | | | | | | | |
| w | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

We set: $Z = \log^2 m$, $\tau = \frac{1}{2} \log m$

Big zero small zero

then, the space:

$$O\left(\frac{m}{\log^2 m} \cdot \log m + \frac{2m}{\log m} \cdot \log(\log^2 m) + \left(\frac{1}{2} \log m\right) \cdot 2^{\frac{1}{2} \log m} \cdot \log\left(\frac{1}{2} \log m\right)\right)$$

$$O\left(\frac{m}{\log m} + \frac{m}{\log m} \log \log m + \underbrace{(\log m) \sqrt{m} (\log \log m)}_{m^{1/2 + \epsilon}}\right) = O(m)$$

($m \log \log m / \log m$)

alias - Four idea to compress

$$L = \begin{array}{|c|c|c|c|c|} \hline & & & \dots & \\ \hline l & l & l & & \\ \hline \end{array}$$

$n \lceil \log_2 \frac{m}{n} \rceil$ bits

Th : $H = \underline{11101100\dots}$

$2m$ bits

Select₁ } + $O(n)$ bits
Select₀ }

Th 2: Elias-Fano for encoding an integer array can be implemented

in $2m + n \lceil \log_2 \frac{m}{n} \rceil + O(n)$ bits and support access in $O(1)$ Time

and NextKey in $O(\log \frac{m}{n})$ time

$$O(m(2 + \lceil \log_2 \frac{m}{n} \rceil + o(1)))$$

Given $B[1, m]$, construct $A[1, m]$ such that $A[i] = \text{position of } i\text{-th 1}$.

$$B = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 = m \\ \hline 0 & 1 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad O(m) \quad m=6$$

$$A = \begin{array}{|c|c|c|c|} \hline 1 & 2 & 2 & 4 = m \\ \hline 2 & 3 & 5 & 6 \\ \hline \end{array} \quad O(n)$$

Th 1 } $m + O(m)$ bits
} $O(1)$ rank/select

We build A putting the positions of 1 occurring in B.

build Th 2 over A \rightarrow Elias-Fano + select over array H,

$m = \max \text{pos in } B$

$$\lceil m(2 + \lceil \log \frac{m}{n} \rceil) + O(1) \rceil \text{ bits}$$

Th 2

Scans $O(1)$
Or $O(\log \frac{m}{n})$ bits

$\text{Select}_1(i) = \text{Access}_{\frac{m}{n}}(i)$ in $O(1)$ time

$\text{Rank}_1(s) = \text{GQ}_{\frac{m}{n}}(s+1) - 1$ in $O(\log \frac{m}{n})$ time

$S = \boxed{\text{B} \text{ O } \times \text{ O } \text{ C } \text{ A } \text{ T } \text{ O } \text{ P } \text{ A } \text{ O } \text{ L } \text{ O } \text{ V } \text{ O } \dots}$

$A = \boxed{1 \ 5 \ 9 \ \dots}$

$B = \boxed{1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 1}$

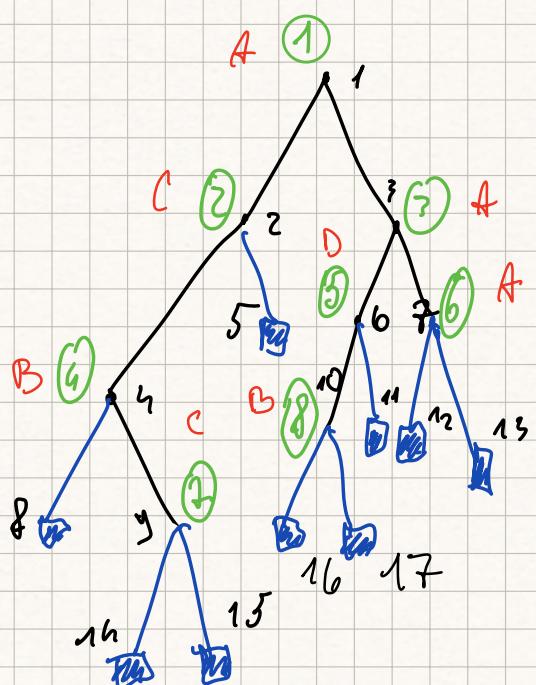
Build Sol₁ compact

- $D + O(D)$ bits
- Retrieve a string i -th: Scan S from position: $\text{Select}_1(B, i)$

Build Sol₂ compressed

- $m \left(2 + \log \frac{D}{m} + O(1) \right)$

Binary tree: Succinct representation of binary trees
(pointers programming)



- ① tree completion ← *empty leafs*
- ② node labeling (breadth-first) ← *add dummy nodes*
- ③ serialization (1 if the position corresponds to an existing node)

| | | | |
|-------|--|---|-------------------|
| $B =$ | $\begin{matrix} 1 & 2 & 4 & 5 & 6 & 2 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \\ 1 & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$ | \Rightarrow^m | $\text{Row } K_1$ |
| | $m=1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8$ | $+ \quad \quad \quad \quad \quad \quad \quad \quad \quad$ | Select_1 |

Total space occupancy: $B + O(|B|)$ $|B| = m = 2m + 1$ bits

$$m(2+O(1))+1 = 2m+1+O(m) = m+O(m) \text{ bits } O(1) \text{ time}$$

only 2 bits per node

$2^m \geq \log_2 2^m \geq 2^m$ bits to represent a tree

$\text{left-child}(x) = \text{Rank}_1(2x)$ skip the fake nodes (leaves)

not circled if $B[2x] = 1$

$\text{right-child}(x) = \text{Rank}_1(2x+1)$ if $B[2x+1] = 1$

$$\text{parent}(x) = \left\lfloor \frac{\text{select}_1(x)}{2} \right\rfloor$$

$L =$

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| A | C | A | B | D | A | C | B |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

Date comprehension:

'50-'60 . Statistical compressors (Huffman, arithmetic, ...)

'70 . Dictionary compressors (Lempel-Ziv family $\leftarrow \begin{smallmatrix} LZ77 \\ LZ78 \end{smallmatrix} \right)$

'80 . Block sorting compressors (Burrows-Wheeler transform ...)

gzip, tarzip

Compressed-file :

| | |
|-----------|------|
| PROGAMBLB | BOOY |
|-----------|------|

\sum, prob

01001...

We discuss this in the following
lecture

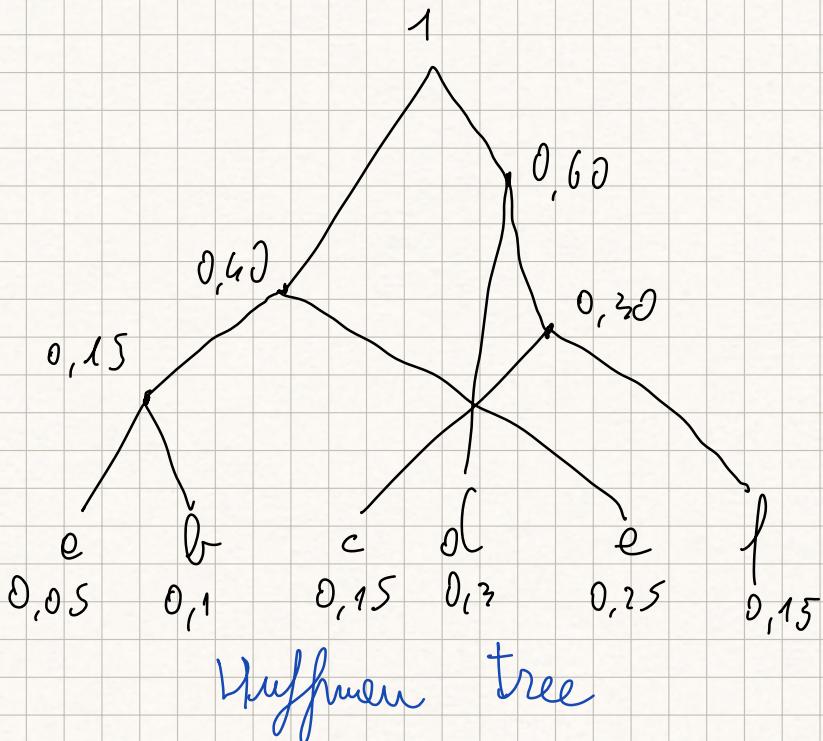
0,3
1
0,15

$$\sum = \{a, b, c, d, e, f\} \text{ sum} = 1$$

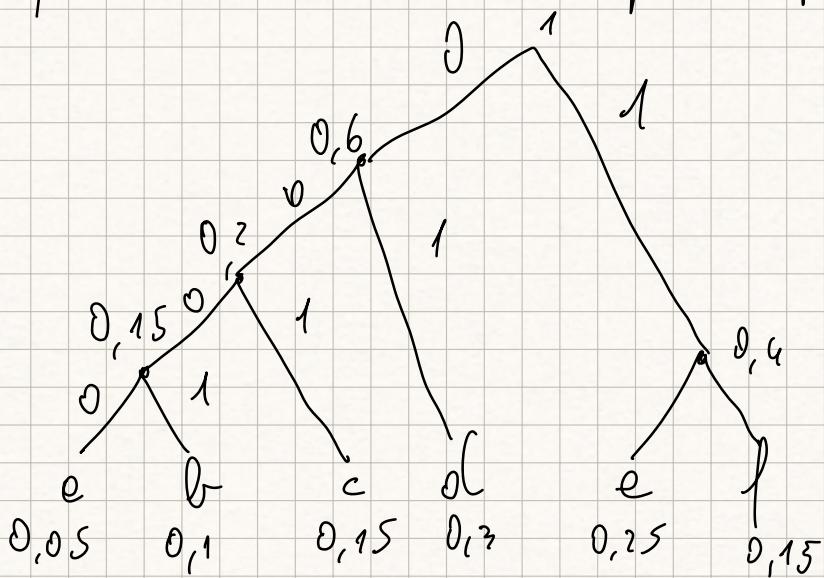
| | | | |
|------|-----|------|------|
| 0,05 | 0,1 | 0,15 | 0,15 |
|------|-----|------|------|

Huffman: ① pick the two nodes with smallest probability
 ② merge them into a parent node

③ trying to prevent the sum of prob.



encoding of e will be 0000, the path from root to leave



The: Huffman code is optimal among all prefix-free code for alphabet

$$\sum$$

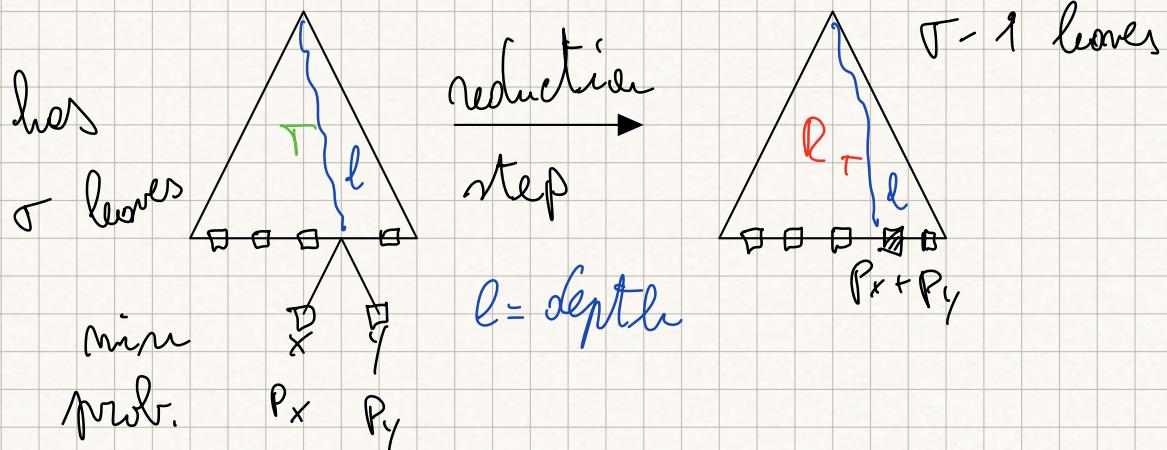
$$\uparrow$$

$$L_H \leq L_C \quad \forall C \text{ prefix-free code}$$

Lemme 1: F = set of binary trees with Σ leaves where avg depth is minimum. If tree $\in F$ such that the two leaves with min prob. are at the largest depth and are children of the same parent.

avg depth of a leaf = avg codeword length $L_c = \sum_{s \in \Sigma} p(s) \cdot |\text{cw}(s)|$

Lemme 2:



L_T = average depth tree T $\sigma^{fx,y}$

$$L_T = \sum_{\sigma} p(\sigma) \cdot l(\sigma) = \sum_{\sigma \neq x,y} p(\sigma) l(\sigma) + p(x)l(x) + p(y)l(y) =$$

$$= \sum_{\sigma \neq x,y} p(\sigma) l(\sigma) + p(x)(l+1) + p(y) \cdot (l+1) =$$

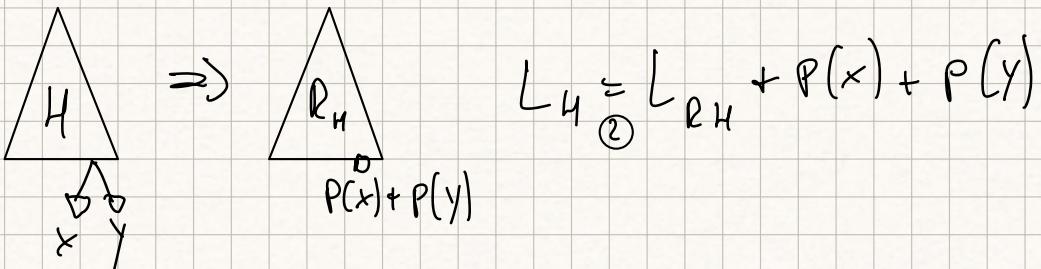
$$= \sum_{\sigma \neq x,y} p(\sigma) \cdot l(\sigma) + [p(x) + p(y)] \cdot (l+1)$$

$$L_R = \sum_{\sigma \neq x,y} p(\sigma) \cdot l(\sigma) + l(p(x) + p(y))$$

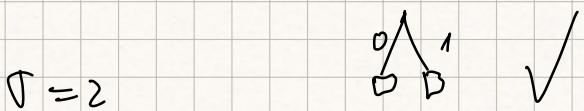
$$L_T - L_R = (\ell+1) \cdot p(x) + (\ell+1) \cdot p(y) - \ell(p(x) + p(y)) = p(x) + p(y)$$

$$L_T = L_R + p(x) + p(y)$$

Direkt Huffman thm:

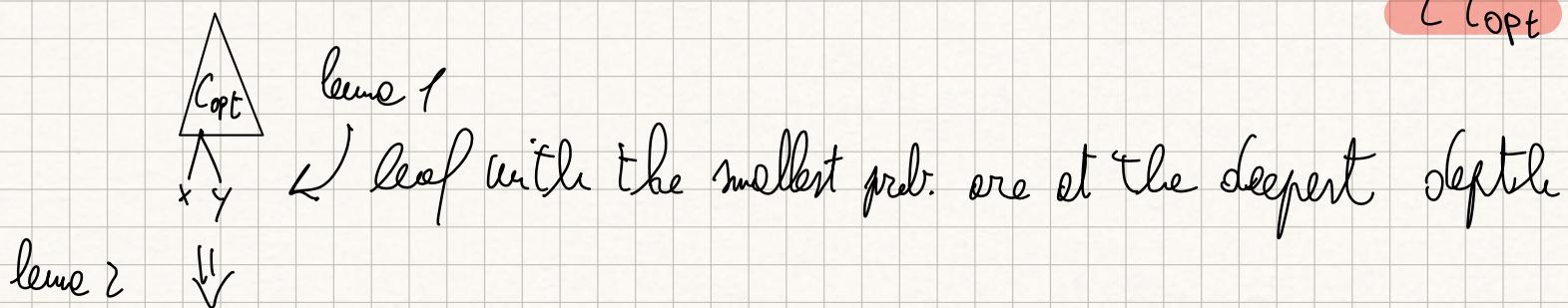


Proof by induction:



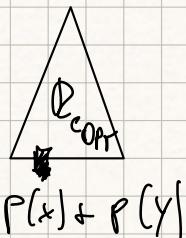
Huffman is optimal for an alphabet of $r-1$ symbols

let C_{opt} an optimal code for r symbols \Rightarrow avg depth is L_{opt}



leaf with the smallest prob. are at the deepest depth

level 2



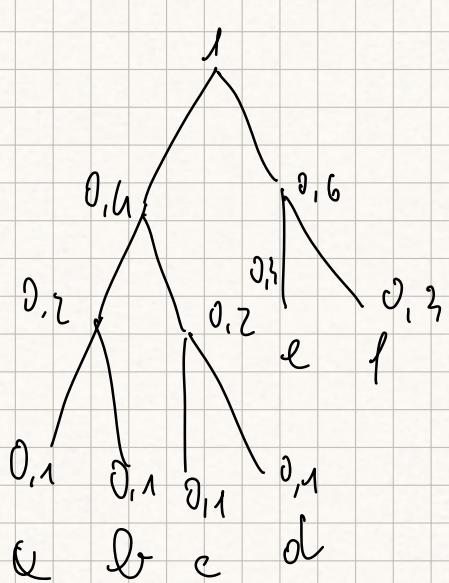
$$L_{H,\sigma} \stackrel{(2)}{=} L_{R_{H,\sigma-1}} + p(x) + p(y), \leq L_{R_{\text{OPT},\sigma-1}} + p(x) + p(y) = L_{\text{OPT},\sigma}$$

optimal for $\sigma-1$
 symbols according
 to the inductive
 hyp.

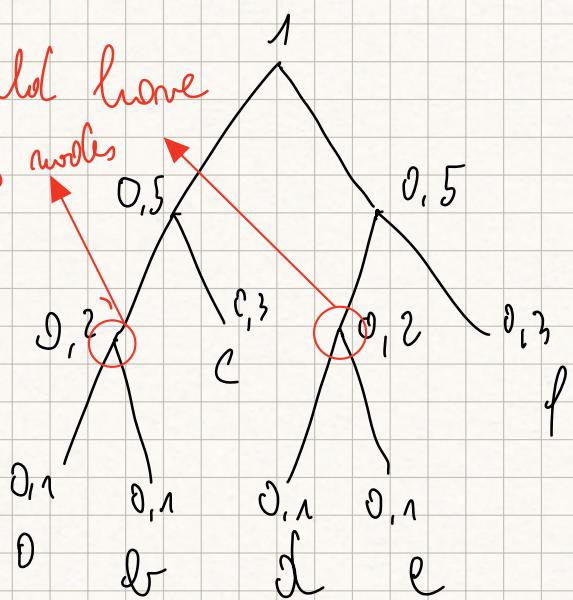
smallest prob.
 over σ symb.

$L_{H,\sigma} \leq L_{\text{OPT},\sigma}$ ■

Huffman tree is shallower than the other



We should have merged this nodes



Huffman tree

Entropy

$$H_{\text{Entropy}} \leq L_{\text{Huff}} < H+1 \text{ bits}$$

Shannon

avg codeword length

every codeword on avg loses < 1 bit w.r.t. entropy

$$0 \leq H \leq \log_2 |\Sigma|$$

if we have a repetitive source of symbols (*skewed distribution*)

HUFFMAN isn't the best choice

$$H \rightarrow 0 \quad L_{\text{Huff}} \geq 1 \text{ bit} \quad \text{skewed distribution}$$

$$H \rightarrow \log_2 |\Sigma| \quad L_{\text{Huff}} \approx H \quad \text{uniform distribution}$$

To reach the entropy we incorporate sets together in order to form the macro symbols

$$\frac{1}{k} + \frac{1}{k} + \frac{1}{k} + \dots = |\Sigma'| = |\Sigma^u|$$

avg cw length

1948

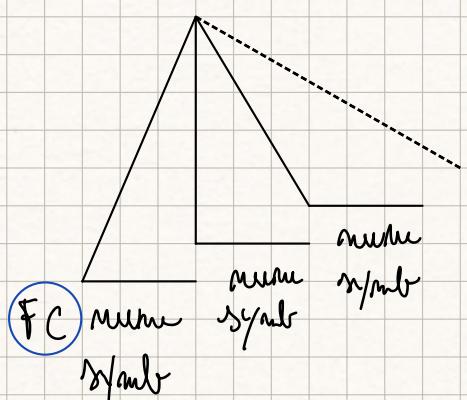
$$H \leq \frac{L_{\text{Huff}}}{K} \quad \text{for a block of length } K \leq H + \frac{1}{K}$$

K

avg number of bits per symbol taken by a cw

Input: Huffman tree

Canonical Huffman tree



First codeword

Given arr [1, maxcw] : array of ints
Symbol [1, maxcw] : array of arrays of symbols

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

1 2 3

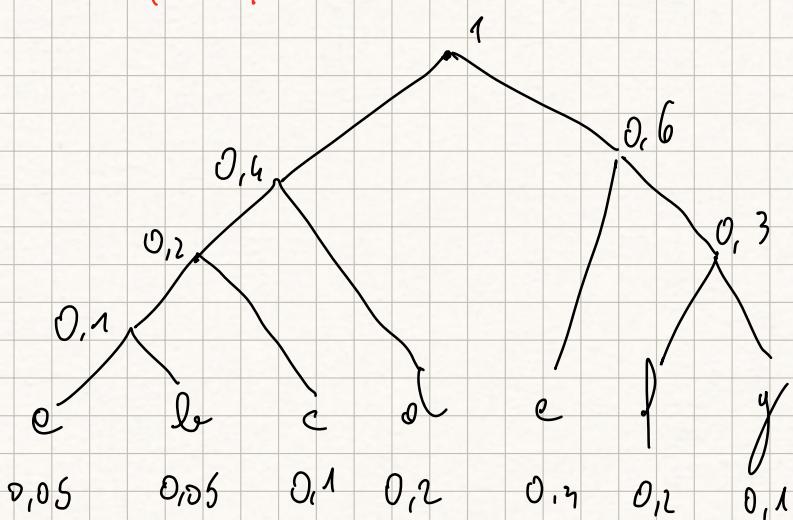
number of symbols per level

$$\text{arr} [1, 3] :$$

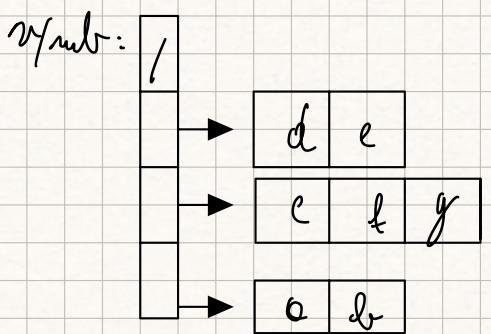
| |
|---|
| / |
| |
| |

| | |
|---|---|
| e | f |
| | |

| | | | |
|---|---|---|---|
| e | b | c | d |
| | | | |



$$\text{num}[1,4] = [0,2,3,2]$$



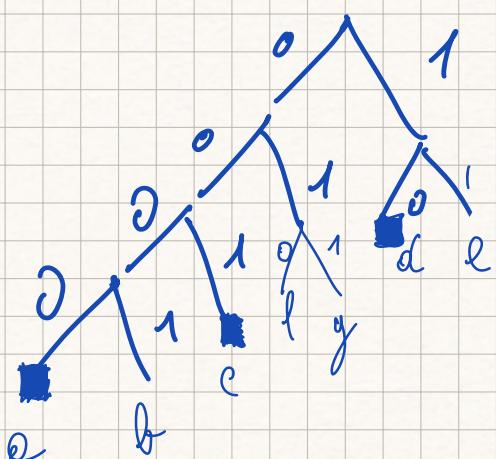
$$FC[\max \text{ level}] = \underline{4}$$

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

$$FC[3] = \frac{FC[4] + \text{num}[4]}{2} = \frac{0+2}{2} = 1$$

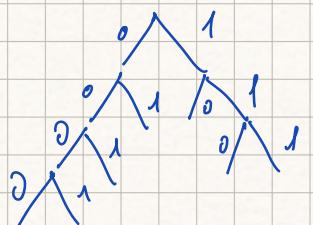
$$FC[2] = \frac{FC[3] + \text{num}[3]}{2} = \frac{1+3}{2} = 2$$

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

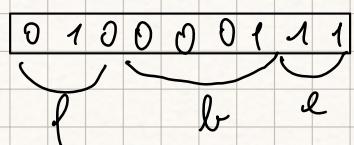


$$FC[1,4] = [2, 2, 1, 0]$$

1 2 3 4



Comprendere:



$$\cdot \left(FC[\text{level symbol}] + \text{offset symbol} \right)_2$$

Pseudocodice:

Decompress ():

v = next_bit(); // prendo il primo bit dello spettro

l = 1; // parto dal primo livello

while ($v < FC[l]$) // finché (v)₁₀ < ($FC[l]$)₁₀ vuol dire che non
sono arrivato al livello in cui è presente
il simbolo

v = 2 * v + next_bit(); // eseguo cambio di base di v da 10 → 10
l++ // incremento livello

return symb [l, V - FC[l]] // decomprimo e ottengo il simbolo

a: $v=0, l=1$

$v=1, l=2$

$v=2, l=3$

symb[3, 2-1]

b: $v=0, l=1$; $v < FC[1] = ?$

$v=2 \cdot 0 + 0, l=2; v=0 < FC[2] = 2$

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

$v=2 \cdot 0 + 1 = 1, l=4; v=1 > FC[4] = 0$
STOP

symb[4, 1-0] = symb[4, 1] = b

c: $v=1, l=1$; $v=1 < FC[1] = 2$

$v=2 \cdot v + \text{next_bit}() = 3, l=2;$

$v=3 > FC[2] = 2$

symb[2, 3-2] = e

$$\text{Compression ratio} = \frac{\text{len comp. file}}{\text{len orig. file}} = \frac{9}{24} \leq 1$$

file of n symbols $\xrightarrow{\text{best}} n$ bits for the comp. file

$$\text{C. R.} = \frac{1 \text{ m bit}}{8 \text{ m bytes}} = \frac{1}{8} \approx 12.5\%$$

Arithmetic Coding:

Converter(x), $x \in [0, 1) \rightarrow b_1 b_2 b_3 \dots$ $b_i \in \{0, 1\}$

repeat until k bits
are emitted

$x = 2 \cdot x$;
if $x < 1$ then output 0;
else {output 1; $x = x - 1$;}

| | | |
|--------------|-------|---|
| $x = 1/3$ | $2/3$ | 0 |
| \downarrow | $4/3$ | 1 |
| .01 | $1/3$ | |

dynamic fraction $x = \frac{J}{2^t}$, $J \in \mathbb{N}$
 $t \in \mathbb{N}$

$x = 00\dots 0$ bin (J)
 t bits

$$\frac{5}{32} = \frac{5}{2^5} = .\underline{0}\underline{0}\underline{1}\underline{0}\underline{1}$$

$\frac{1}{8}$ $\frac{1}{32}$

Lemma: given $x \in (0, 1)$, if we truncate x to the first d bits

\hat{x} it holds $0 \leq x - \hat{x} \leq 2^{-d}$ error

$$x = .b_1 b_2 \dots b_d | b_{d+1} b_{d+2} \dots \quad \left\{ \begin{array}{l} 0 \leq x - \hat{x} \leq 2^{-d} \\ \hat{x} = .b_1 b_2 \dots b_d | 0 \quad 0 \dots \end{array} \right.$$

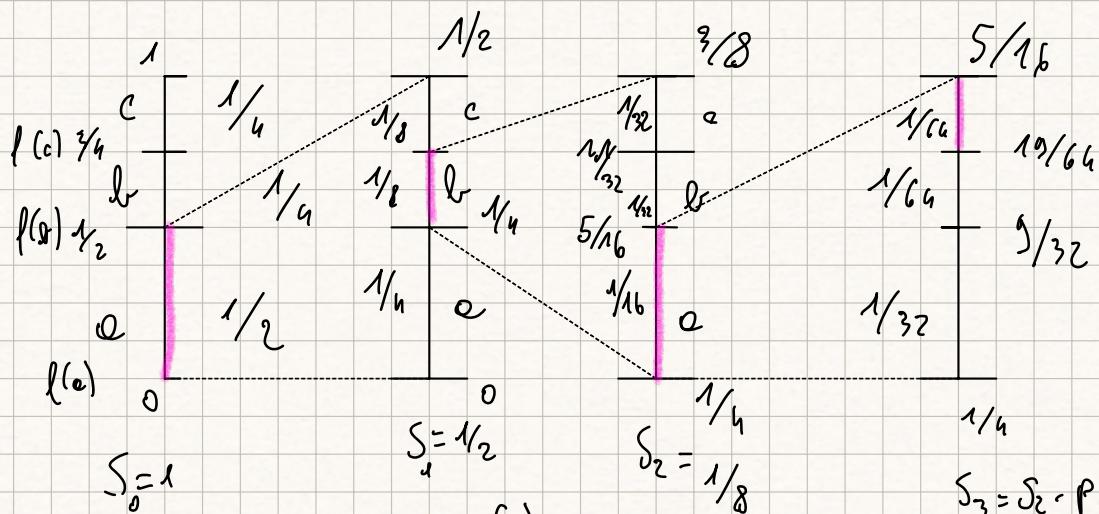
$$\begin{aligned} x - \hat{x} &= .0 \dots 0 b_{d+1} b_{d+2} \leq \frac{00 \dots 0 111 \dots}{2^d} = \sum_{i=1}^{+\infty} 1 \cdot 2^{-(d+i)} = \\ &= 2^{-d} \cdot \underbrace{\sum_{i=1}^{+\infty} 2^{-i}}_1 = 2^{-d} \blacksquare \end{aligned}$$

$$T = e \ b \ c \quad \text{semi-static model} \quad P(e) = \frac{2}{6} = \frac{1}{3} \quad P(b) = P(c) = \frac{1}{6}$$

$$S_0 = 1 \quad S_i = S_{i-1} \cdot P[T[i]]$$

$$L_0 = 0 \quad L_i = L_{i-1} + S_{i-1} \cdot P[T[i]]$$

$$\begin{aligned} f(e) &= 0 \\ f(b) &= P(b) = \frac{1}{2} \\ f(c) &= P(b) + P(c) = \frac{3}{6} \\ f(\Gamma) &= \sum_{c \in \Gamma} P(c) \end{aligned}$$



$$T = e \ b \ c \rightarrow \left[\frac{10}{64}, \frac{5}{16} \right) \quad L_1 = \frac{1}{16}$$

$$J_m = \frac{1}{64} \quad L_m = \frac{10}{64}$$

$$S_3 = S_2 \cdot P(e) = \frac{1}{16} \quad L_3 = \frac{1}{16}$$

$$= S_1 \cdot \frac{1}{4} P(b)$$

$$= S_0 \cdot P(a)$$

$$L_1 = 0$$

$$S_0 = 1$$

$\{$
 - $y = \text{pick a number inside}$
 - $n = \text{length of } T$

$$S_0 = 1 \quad | \quad l = 2 \quad |$$

for $i=1$ to n

$$S_i = S_{i-1} \cdot P(T[i])$$

probabilità del simbolo i -th
delle sequenze

$$l_i = l_{i-1} + S_{i-1} \cdot f(T[i])$$

dimensione intervallo i -th simbolo
delle sequenze

return $\text{comacter}(l_m + \frac{S_m}{2})$ unit $\lceil \log_2 \frac{2}{S_m} \rceil$ bits

Truncating:

$$\text{let's assume } d = \lceil \log_2 \frac{2}{S_m} \rceil \xrightarrow{\text{base}} 2^{-d} = \frac{1}{2^{\lceil \log_2 \frac{2}{S_m} \rceil}} \leq \frac{1}{2^{\log_2 \frac{2}{S_m}}} = \frac{1}{2} = \frac{S_m}{2}$$

$$x = \frac{13}{64} + \frac{1}{128} = \frac{39}{128} \text{ number to send for compression}$$

$$x = l_m + \frac{S_m}{2}$$

$$d = \lceil \log_2 \frac{2}{1/64} \rceil = \lceil \log_2 128 \rceil = 7 \cdot \underbrace{0100111}_2 \text{ bits to encode } \hat{x}$$

Th: the number of bits emitted by arithmetic is at most $2 + n \cdot H$

$$\lceil \log_2 \frac{2}{S_m} \rceil = \lceil 1 - \log_2 S_m \rceil < 2 - \log_2 S_m = 2 - \log_2 \prod_{i=1}^n P(T[i]) =$$

$$= 2 - \sum_{i=1}^n \log_2 P(T[i]) = 2 - n \sum_{\sigma} \frac{\text{occ}(\sigma)}{m} \cdot \log_2 P(\sigma) = 2 + n \sum_{\sigma} P(\sigma) \log_2 \frac{1}{P(\sigma)} =$$

$$= 2 + n \sum_{\sigma} P(\sigma) \log_2 \frac{1}{P(\sigma)} = n \cdot H + 2 \text{ bits} \blacksquare$$

$$\# \text{ bits per symbol} \frac{n \cdot H + 2}{n} = H + \frac{2}{n}$$

Huff. compressed len $\leq n(H+1) = nH + n$ bits

Arith. compressed len $\leq nH + 2$ bits

Huff. per symbol $\leq H+1$

Arith. per symbol $\leq H + \frac{2}{n}$

Lempel-Ziv '77 '78 (dictionary based compression) \rightarrow gzip, 7z, brotly

- dictionary (\rightarrow all substrings, \rightarrow some substrings) (77) (78)
- updating procedure for the dictionary

transform
the text

peers
pairing strategy

$T \rightarrow \langle d, l, \text{len}+1 \rangle$

$T = A A C A A C A B C A A A A A C$

dictionary = set of all substrings starting before

$\langle 6, 3, A \rangle$ $\langle 1, n, c \rangle$

string for d

string for l

string for len

$S_d = 0 \ 1 \ 3 \ 6 \ 4 \ }$

$S_l = 0 \ 1 \ 4 \ 3 \ 4 \ }$

$S_{\text{len}} = A \ C \ B \ A \ C \ }$

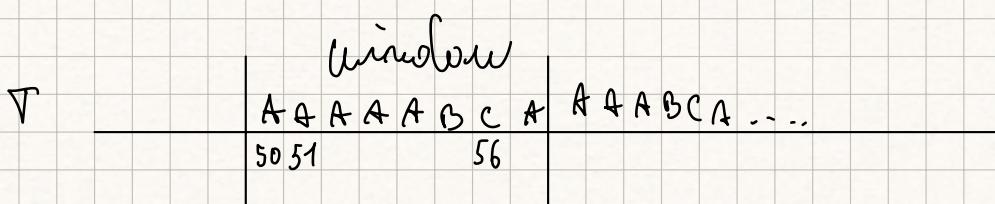
cur-d

Huffman / Arithmetic

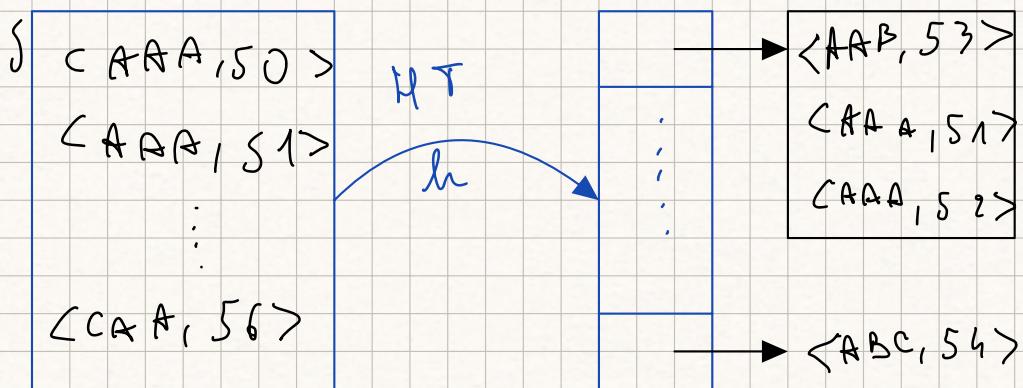
$A \rightarrow C \rightarrow D$ | $C \rightarrow D \rightarrow C \rightarrow D \rightarrow C \rightarrow E$
 cur ↑ (C, D, E)

for ($i=0$, $i < len$, $i++$)

$$\text{out}[cur+i] = \text{out}[cur - \delta + i]$$



large $w \rightarrow$ better compression



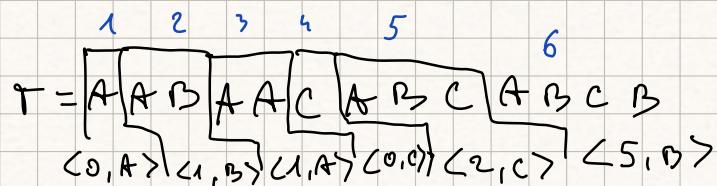
qtip implements configuration that change the window of compression

window = 1 100kb fast compression speed

poor compression ratio

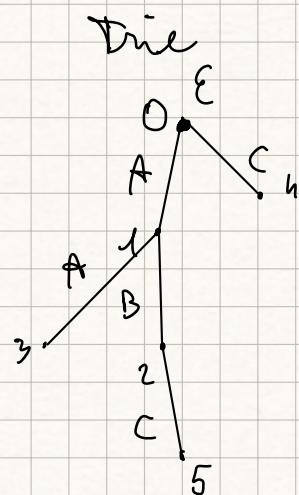
window = 8 900kb medium compression speed
 better compression ratio

L278 : we construct a tries data structure



→ find the longest copy in the dictionary $\equiv S$

- init ID of the longest copy found so far next char (pair)
- update the dictionary, adding S_c with new ID



Text formed by A^m char:

L272

② $\langle 0, 0, A \rangle; \langle 1, m-1, 0 \rangle;$

L278

\sqrt{m}

$$1 + 2 + 3 \dots + \sqrt{m} = m$$

$$\frac{\sqrt{m}(\sqrt{m}+1)}{2} = m$$

Brotli, 25+10

ANS

Asymmetric
Numerical
System

General
Histogram

$L = \{P, S, M, P, S, S, I\}$

bzip

$\hookrightarrow \langle \{P, S, M, P, S, S, I\}, 6 \rangle$
L' presable

decompressing ① BWT

local bengenerity ② Move to front

into Wheeler ③ Run length encoding (RLE)

Move-to-Front coding (symbols \rightarrow numbers)
transforms

wheeler's code

input: (string, list of symbols)

$$\begin{cases} L' = \{P, S, M, P, S, S, I\} \rightarrow \begin{array}{ccccccccc} 1 & 2 & 3 & 4 & 1 & 4 & 3 & 4 & 1 \\ 10 & 10 & 10 & 10 & 10 & 10 & 10 & 10 & 10 \end{array} \geq 2 \text{ at least two bit} \\ L = \{1, 2, 3, 4\} \text{ presable} \quad \text{SIP } \not\equiv \text{ m} \quad \text{drop the first bit} \end{cases}$$

symbols \rightarrow probabilities
 $\begin{array}{ccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \downarrow & \downarrow \\ 10 & 10 & 10 & 10 & 10 & 10 & 10 & 10 & 10 \end{array}$
 $P = \frac{1}{12}, \frac{1}{12}, \frac{1}{12}, \frac{1}{12}, \frac{1}{12}, \frac{1}{12}, \frac{1}{12}, \frac{1}{12}, \frac{1}{12}$
 $\sum = \{0, 1, 2, 3, 4\}$
 Then apply Huffman / Calkin-Wilf