

SORTING STRINGS

- strings form an alphabet Σ
 - $|\Sigma| = \infty$
 - $S[1, n]$ array of pointers to strings \hookrightarrow variable length
 - N : total length of strings
 - M : # strings
 - AVG length $L = \frac{N}{M}$
- $\frac{1}{\downarrow}$ $\frac{1}{\downarrow}$
 pointers strings are
 are spread
- $\Rightarrow \text{QSort}(\text{arrays}, M, \text{string objects},$
 $\swarrow \quad \text{Comparator-function})$
- This takes \Rightarrow comparisons \circ AVG length of strings
- $$\Rightarrow O(L \cdot m \log n) =$$
- $$= O\left(\frac{N}{M} \cdot m \log n\right) =$$
- $$= O(N \log n)$$

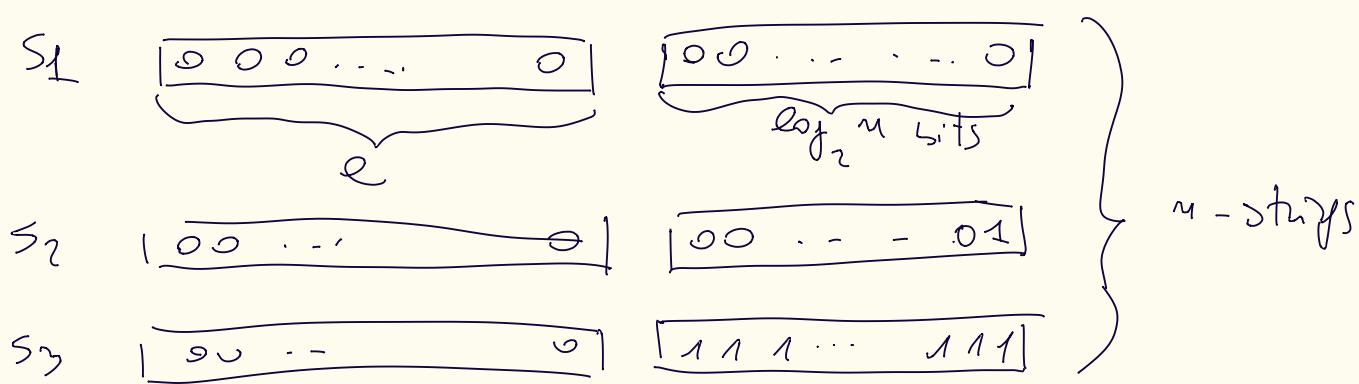
$\frac{1}{\downarrow}$
NOT OPTIMAL

LOWER BOUND

- ① to sort n strings \Rightarrow SORT 1st CHARACTER $\Rightarrow \Omega(n \log n)$
 - ② $d = \sum_s d_s \Rightarrow \Omega(d)$
 distinguish prefix of a string $\log n$ times
- $\Rightarrow d = \sum_s d_s$
 \Rightarrow lower bound $\Omega(n \log n + d)$ Time

• IMPACT OF lower bound of st�ys to QS (comparison-based algorithm)

- given set of st�ys:



$$d_i = l + \log_2 n$$

1
distinguish
prefix

$$\Rightarrow d = \sum_{i=1}^n (l + \log_2 n) = l \cdot n + \log_2 n$$

↓
overall
string

lower bound

$$\Omega(f(n \log n + d))$$

$$QSORT = \left(\frac{N}{m} \cdot l \cdot m \log n \right) \quad \text{in this case } l = l + \log n$$

$$\Rightarrow QS = l \cdot m \log n + m \log^2 n$$

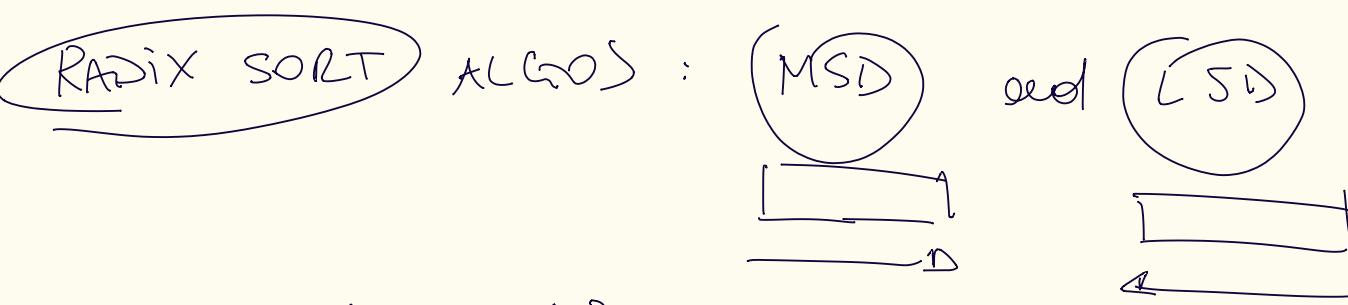
this case

$$\begin{aligned} & QSORT \left(l \cdot m \log n \right), \\ & = (l + \log n) (m \log n) \end{aligned}$$

$$O(l m \log n + m \log^2 n) \text{ UPPERN BOUND}$$

QSORT IS FAR FROM lower BOUND

- QSORT IS AT LEAST $\sqrt{\log n}$ Factor from lower bound

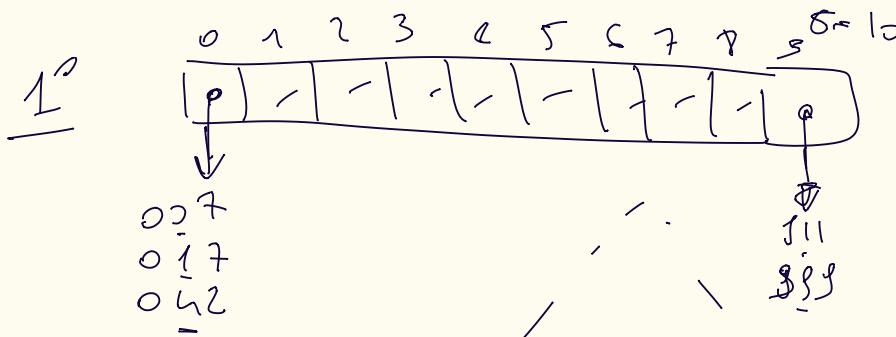


• Stays or digits

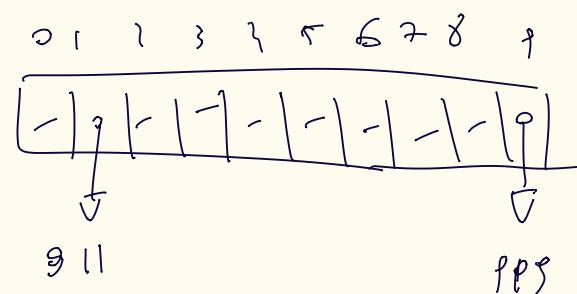
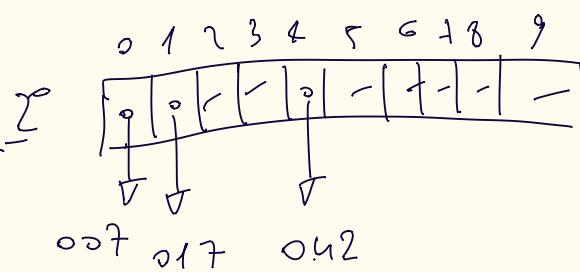
MSD

$$S = \{007, 017, 042, 911, 999\}$$

→ sum according to 1st symbol



→ sum according second digit



- step: every socket 1 stays

VISIT

SORTED SONGS

SIZE

stay s we create ds nodes

$$\# arrays = \# nodes \leq d = \sum_s d_s$$

↓ not efficient

$$\text{SPACE} = O(d \cdot \delta)$$

↓ efficient

$$\text{TIME } O(d \cdot \delta)$$

but usually δ is constant $\Rightarrow O(d)$ time

\Rightarrow \forall node I'm allocating
the position (order of size δ)

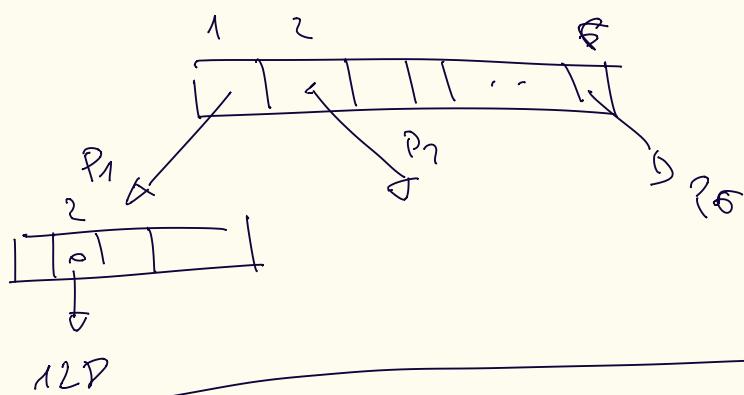
$\Omega(n \log n + d)$ \Rightarrow why this is comparison-based lower bound

\Rightarrow here MSD is $\mathcal{O}(d)$ \Rightarrow this is radix-based

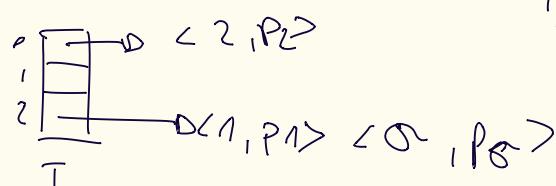
cool for scanning $\mathcal{O}(p)$ time \Rightarrow look the pattern
not a comparison based
can implement MSD with hash with chaining

- pair: < digit, pointers > \rightarrow inserted in hash table
 $\begin{array}{c} \text{digit} \\ \downarrow \\ \text{key} \end{array}$ $\begin{array}{c} \text{pointers} \\ \downarrow \\ \text{satellite data} \end{array}$

- create an hash table of every node whose size is proportional to # of distinct pairs (# edges)



\Rightarrow note:



$\begin{array}{c} h \\ \langle 1, p_1 \rangle \\ \vdots \\ \langle 12, p_{12} \rangle \end{array}$ $\begin{array}{c} h \\ \langle 2, p_1 \rangle \\ \vdots \\ \langle 12, p_{12} \rangle \end{array}$ $\begin{array}{c} h \\ \langle 6, p_1 \rangle \\ \vdots \\ \langle 12, p_{12} \rangle \end{array}$
 $\approx 612 \text{ bytes}$

$$|T| = \Theta(\# \text{keys})$$

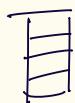
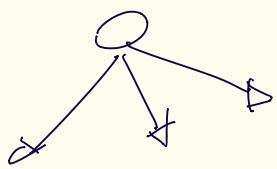
\Rightarrow cost of search on T is $\mathcal{O}(1)$ per op

$$\frac{M}{m} \times \frac{3}{3}$$

$$\Rightarrow \text{Space} = \Theta(\# \text{edges}) = \Theta(\text{keys})$$

old have
on implementation
with hash
 \Rightarrow but digest are not used

MSD with hash table



- on hash table + node

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

1) Build tree $O(1)$ insertion of every strip (hit)

and tree to do + strip

$$\sum_s O(1) * d_s \Rightarrow O(d) \text{ AVG}$$

↑ for overall strips

2) Sorting hash tables (edges) \Rightarrow Q SORT on the pairs

$$d_u \cdot g_m$$

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

$$O(\underbrace{\# \text{ edges}}_{e_v} \cdot \lg \# \text{ edges})$$

$$\sum v_i \log_2 v_i \leq \sum v_i \log \sigma \Rightarrow O(v_i \lg v_i)$$

$$= \log_2 \sigma \left(\sum_u v_u \right) = O(\lg \sigma \cdot d) =$$

$$\Rightarrow O(d \lg \sigma) \text{ time}$$

SPACE $O(d)$

LOWER BOUND FOR SORTING STRINGS

⇒ Sorting according to 1st digit

$$\Omega(n \log_2 n)$$

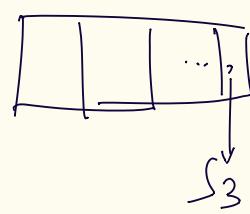
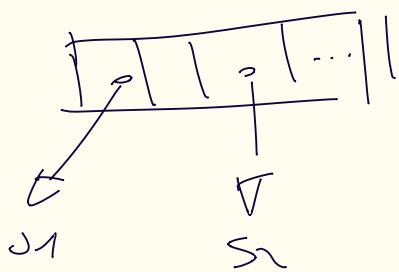
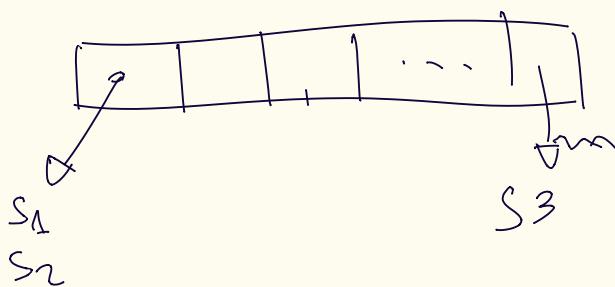
⇒ Consider distinguish prefixes

$$\sum_s f_s = d \quad \mathcal{O}(d)$$

⇒ lower bound

$$\lceil \Omega(n \log n + d) \rceil$$

→ Consider MSD radix sort



$\sigma = |\Sigma|$
 - given the strings
 ⇒ ~~also~~ create ~~and~~ array of
~~σ - positions~~
 ⇒ distribute the
 strings acording
 to the 1st char,
 the 2nd and so on
 until you have 1 string
 per bucket of the
 array

⇒ requires
 $= O(d \cdot \sigma)$ time and complexity

⇒ beats lower bound out the lower bound is referred to a comparison-based algos!

\Rightarrow Improve MSD with hash table

\Rightarrow \forall node \Rightarrow store in hash table

\Rightarrow edges \hookrightarrow # keys

Create the tree. $O(1) \cdot \sum_s d_s$

$\Rightarrow O(d)$ on AVG TIME

\Rightarrow Have to sort hash tables

$O(n \log n) =$

$\approx O(\# \text{edges} \log \# \text{edges}) =$

$= O(l_u \log l_u)$

let $b \in \# \text{edges} =$
 $= e_u$

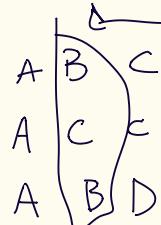
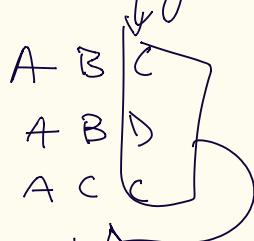
$\Rightarrow \sum_u e_u \log e_u \leq \underbrace{\sum_u e_u \log \sigma}_{O(d)} =$

$\Rightarrow O(d \cdot \log \sigma)$

LSD Radix sort

- STRINGS have same length L

- SORT according to least significant by 1st, 2nd, ... and so on
 (counting sort)



...

Final sort GSTS

$$O(n + \text{max. integer size})$$

\Rightarrow here strings have same length

$$\Rightarrow O(L(n + \epsilon)) \quad \text{time} \quad \text{space } O(N)$$

\Rightarrow can capture more digits per tree

(2)

$$\Rightarrow O\left(\frac{L}{2} \cdot (n + \epsilon^2)\right)$$

(3)

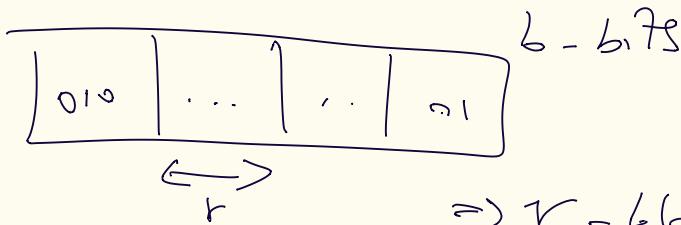
$$\Rightarrow O\left(\frac{L}{3} \cdot (n + \epsilon^3)\right)$$

digit's f

tree ↑

\Rightarrow work on

BITS instead of digits



$\Rightarrow r - 6 \text{ blocks } L - \text{bits}$

$S \leq 2$
 $\Rightarrow \Theta\left(\frac{b}{r} \cdot (m + 2^r)\right)$ want to find R minimize the future

$$2^r \geq m \Rightarrow \frac{b \cdot 2^r}{r}$$

$$r \geq \lg_2 n$$

$\Rightarrow S$ minimized when $r = \Theta(\lg_2 n)$

$$\Theta\left(\frac{b \cdot 2^{\lg_2 n}}{\lg_2 n}\right) = \frac{b \overset{\text{length}}{\cancel{n}}}{\lg_2 n} = \frac{N}{\lg_2 n}$$

MULTIKEY QUICKSORT

- sorting strings & variable length
- R : set of strings
- matches $\Omega(d + m \lg n)$

Algorithm 7.1 MULTIKEYQS(R, i)

```

1: if  $|R| \leq 1$  then
2:   return  $R$ ;
3: else
4:   choose a pivot-string  $p \in R$ ;
5:    $R_< = \{s \in R \mid s[i] < p[i]\}$ ;
6:    $R_= = \{s \in R \mid s[i] = p[i]\}$ ;
7:    $R_> = \{s \in R \mid s[i] > p[i]\}$ ;
8:    $A = \text{MULTIKEYQS}(R_<, i)$ ;
9:    $B = \text{MULTIKEYQS}(R_=, i + 1)$ ;
10:   $C = \text{MULTIKEYQS}(R_>, i)$ ;
11:  return the concatenated sequence  $A, B, C$ ;
12: end if

```

Proof: Consider invariant: all strings in A ($R_<$)
 ↓
 For no
 shares $i - 1$ characters prefix

string $s \in R$:
 $\rightarrow s \in R_< \vee s \in R_> \Rightarrow$ #strings longer
 → occurs
 \lg_2^m times

$\rightarrow s \in R_= \Rightarrow i++$

no other strings share prefix $i \leq L$

thus situation occurs

if pivot is good

$$\Rightarrow \sum_s d_s + \lg_2^m = \sum_s d_s + \overbrace{\sum \lg^m}^{we have m digits, whenever} = d + m \lg m$$

$$= d + m \lg m$$

SORTING USING lower-bound

TAKE ACCOUNT

comparisons and dist. prefix $\sum_s d_s = d$

$\Omega(n \log n)$

$\Omega(d)$

$\Rightarrow \Omega(d + n \log n)$

\hookrightarrow lower bound for comparisons - splitting

(MSB)

$$\sigma = |\Sigma|$$

L = Any length of the strings

$$L = \frac{N}{M} \rightarrow \text{length of strings}$$

$$M \rightarrow \# \text{ strings}$$

- distribute strings

or covering Σ

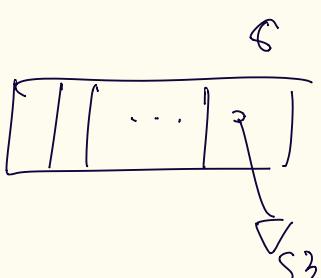
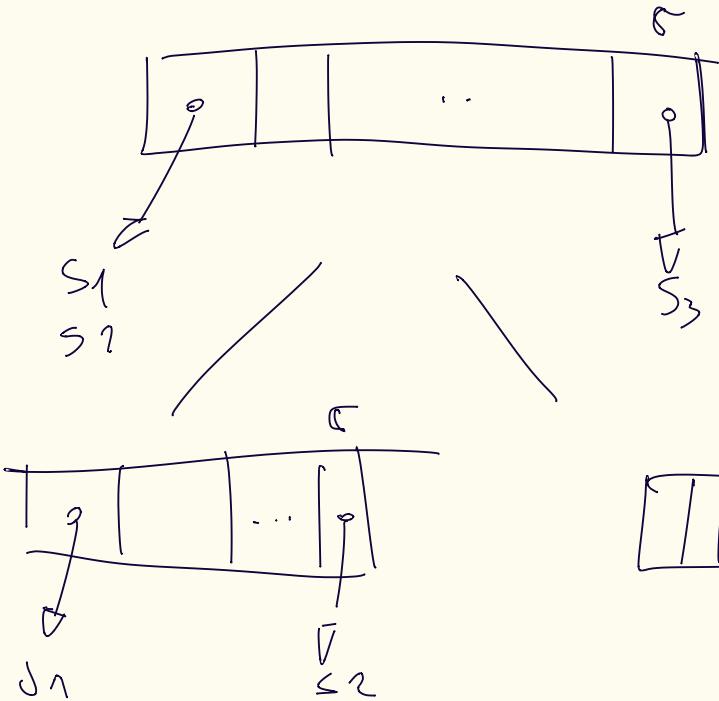
initial character

of the level

\rightarrow next string

$\Rightarrow \text{next string}$

\dots
 \Rightarrow until there
1 string per bullet



\Rightarrow Consider nodes & levels (for comparisons)

$$|\text{array}| = \sigma$$

$$\Rightarrow O(d \cdot \sigma)$$

\Rightarrow I ~~suppose~~ so many levels

as $\# d$

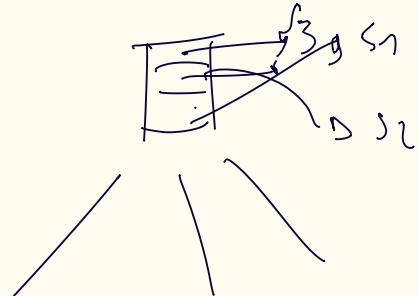
$$\Rightarrow \sum_s d_s = d$$

Improve with hash table \Rightarrow tree

- every node has an hash table

size of HT = ∞

$\Rightarrow h(< \text{string}, \text{id string})$



#edges = $\log n$

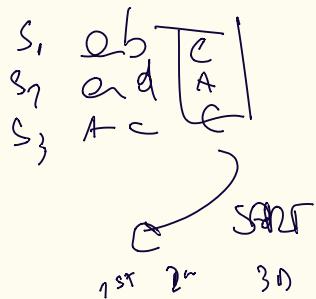
$$O(\underbrace{\# \text{edges}}_{\text{en}} \lg \underbrace{\# \text{edges}}_{\text{en}}) = O(\text{en} \lg \text{en}) =$$

$$\Rightarrow O\left(\sum_s \text{en} \lg \text{en}\right) \leq O\left(\sum_s \text{en} \cdot \lg \text{en}\right) =$$

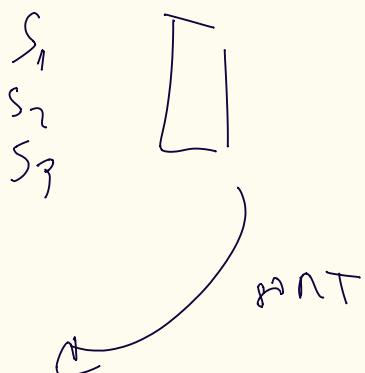
$$\Rightarrow O(d \lg c) \quad (\checkmark)$$

- LSD

- Counting sort $\rightarrow (n + \text{MAX integer size})$



to last significant digit



$$O(L(m + \sigma))$$

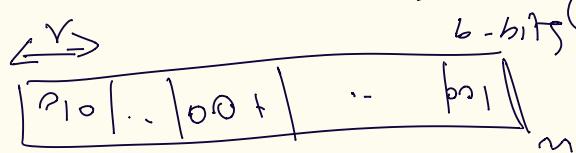
comparing 2 digits at time

$$O\left(\frac{L}{2}(m + \sigma^2)\right)$$

$$\stackrel{\text{step}}{\downarrow} 3 O\left(\frac{L}{3}(m + \sigma^3)\right)$$

$$\stackrel{\text{step}}{\downarrow} 4$$

Consider B.TS instead of digits



Graph n-BTS on
MM

$$O\left(\frac{b}{r}(n + \sigma^n)\right) \quad \sigma = 2$$

$$O\left(\frac{b}{r}(n \cdot 2^r)\right)$$

want a fast min-max factor

$$2^r > n$$

$$r > \log_2 n$$

$$\Rightarrow \left(\frac{b \cdot 2^r}{r} \right)$$

$2^r \leq n$ n^r in cost

$$\Rightarrow \text{Set } r = \Theta(\log_2 n)$$

$$\Rightarrow \Delta \frac{b \cdot 2^r}{r} = \frac{b \cdot 2^{\log_2 n}}{\log_2 n} = \frac{b \cdot n}{\log_2 n} = \frac{L}{\log_2 n}$$

~~MSD~~ = $\left(\frac{L}{\log_2 n} \right)$

MSD ($d \cdot \sigma$)

MSD with hash (trie) = $(d \cdot \log \sigma)$

\sum = MKW-Merkle = $\Sigma (d + m \log n)$

Mult-key \textcircled{Q}

— Matches lower bound

— Strings of m lengths

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

— R : set of strings

if $|R| < 1$ return R

else choose a root prob p in R (p is a string)

$$R_L = \{s \in R : s[i] < p[i]\}$$

$$R_{} = \{s \in R : s[i] = p[i]\}$$

$$R_R = \{s \in R : s[i] > p[i]\}$$

$$A = \text{MKQ}(R_L, i)$$

$$B = \text{MKQ}(R_-, i+1)$$

$$C = \text{MKQ}(R_R, i)$$

Proof invariant: strings in A share $i-1$ character

$$s \in R \xrightarrow{s \in} R_L \vee R_R \Rightarrow \underbrace{\text{occurs}}_{\text{occurs}} \text{lg}_2 m \text{ times} \quad \begin{matrix} \text{-prefix} \\ \text{by } m \text{ times} \end{matrix}$$

$$\xrightarrow{s \in} R_- \Rightarrow i++$$

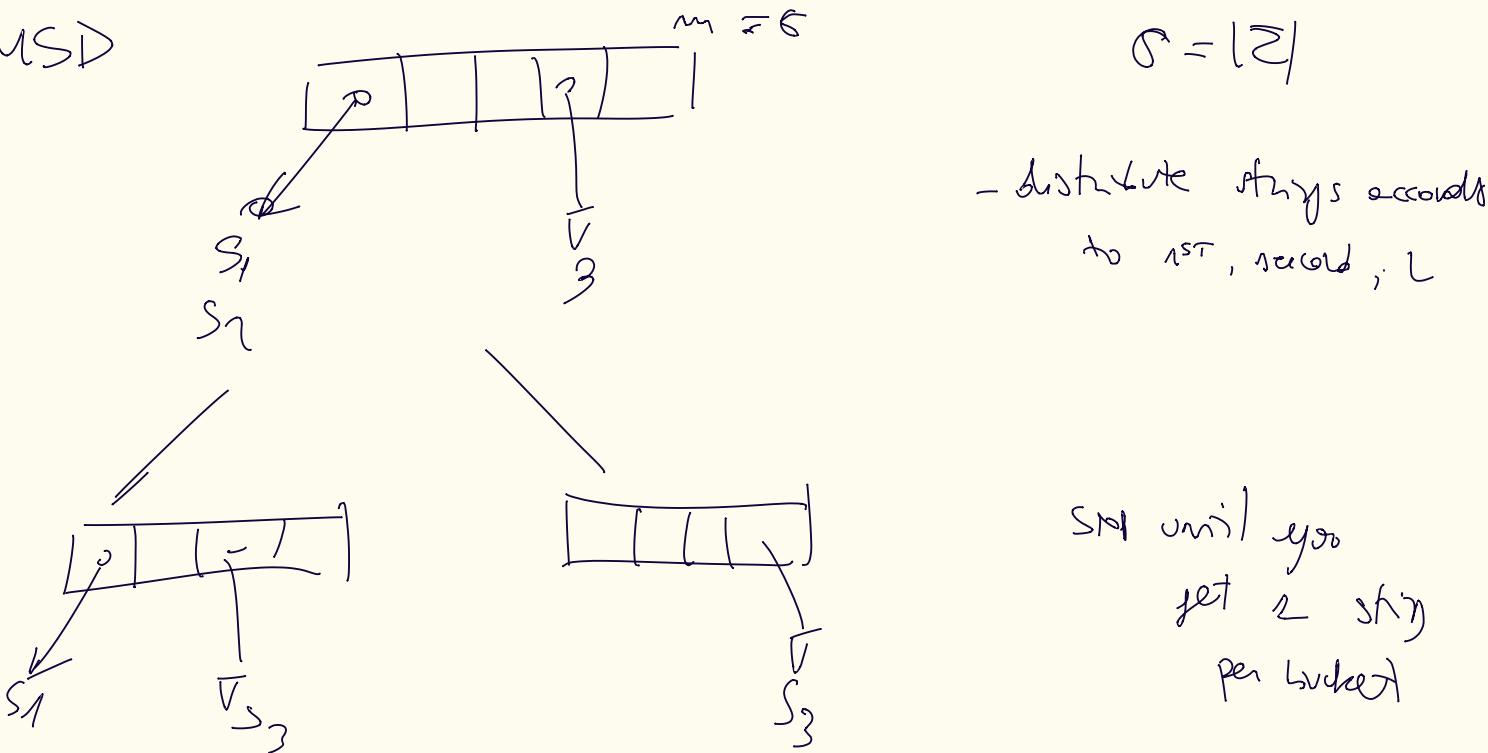
until

$$\xrightarrow{\text{occurs}} \text{when } i = ds \quad i \leq L$$

$$\Rightarrow \sum_s ds + \sum_s \ell f_i M = d(d + n \log n)$$

\downarrow
stays in n

MSD



Read →

=> sorted order

cost: I've array of σ positions
at every level and I've to account for each
the number of distinguish prefixes

$$\Rightarrow d \left(\sum_s d_s \cdot \sigma \right) = O(d \cdot \sigma)$$

\Rightarrow problem some spaces of array are null

\Rightarrow hash \Rightarrow every node has an hash table

(hash or per string, s_1 ,
 s_2 ,
 s_m)

ed

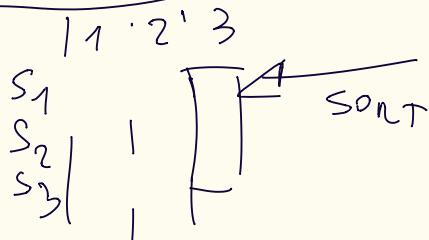
edges:

(for where I make a comparison)

$$O(\# \text{edges} \log \# \text{edges})$$

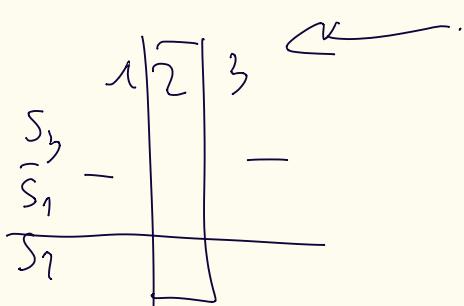
$$\begin{aligned}
 O(\text{eu} \lg \text{eu}) &= O\left(\sum_s d_s \cdot \lg_2 \text{eu}\right) = \\
 &= O(d \cdot \lg_2 \text{eu}) \xrightarrow{\sigma} \\
 &= O(d \lg_2 \sigma)
 \end{aligned}$$

LSD



Counting sort

(n + max representation)



stays have fixed length

and so on

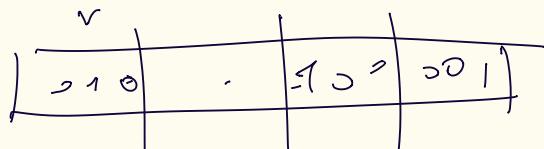
Counting sort

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

can compute τ digits per time

$$O\left(\frac{\tau}{2}(n + \sigma^2)\right)$$

idea: consider bits instead of digits



I've ~~two~~ ~~four~~ digits to count

and I've b-bit

$$\mathcal{O}\left(\frac{b}{r} \cdot (n + 2^r)\right) = \text{but } r \text{ with bits} = 2$$

$$= \mathcal{O}\left(\frac{b}{r} \cdot (n + 2^r)\right)$$

Want r that minimize the factor

$$\Rightarrow r \text{ uses } 2^r > n$$

$$2^r \leq n \text{ not meaningful}$$

$$\hookrightarrow 2^r > n \Rightarrow \frac{b \cdot 2^r}{r} =$$

$$r > \lg_2 n \quad = \frac{b \cdot 2^{\lg_2 n}}{\lg_2 n} =$$

$$= \frac{b \cdot n^L}{\lg_2 n} = \mathcal{O}\left(\frac{n^L}{\lg_2 n}\right)$$

Multilevel QS(R, i)

$R = \text{set of strings}$

if $|R| \leq 1$ RETURN R

else

$p = \text{choose a "good" part from } S, p \text{ is a string}$

$$R_L = \{s \in R : S[i] < p[i]\}$$

$$R_+ = \{s \in R : S[i] = p[i]\}$$

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

$$A = MKQ(R_L, i)$$

$$B = MKQ(R_+, i+1)$$

$$C = MKQ(R_>, i)$$

return A, B, C where $,$ is concatenation

Favoutit: strings in A share $i+1$ chars ~~of~~ suffix

\Rightarrow case: strings in R_L or $R_>$

\Rightarrow occurs by n time

use smy in $R_+ :$

\Rightarrow I've $i+1$

$\Rightarrow i \leq L$

\Rightarrow occurs $\sum_s d_s$ time

$$\begin{aligned} \text{complexity} &\leq \sum_s (d_s + \lg_2 n) = \sum_s d_s + \sum_s \lg_2 n = \\ &= O(d_s + n \lg n) \end{aligned}$$

Multi-key FS

$R = \text{set of keys}$

$\text{MKQ}(R, i)$

if $|R| \leq 1$

return $R;$

else

$p = \text{choose a good pivot } (p \in R)$

$R_L = \{s \in R : s[i] < p[i]\}$

$R_+ = \{s \in R : s[i] = p[i]\}$

$R_R = \{s \in R : s[i] > p[i]\}$

$A = \text{MKQ}(R_L, i)$

$B = \text{MKQ}(R_+, i+1)$

$C = \text{MKQ}(R_R, i)$

$L = \text{left of key}$

Complexity

Invariant: in A stay shorter $i-1$ chars

CASES: 1) stays in $R_L \vee R_R$

\Rightarrow occurs $\log n$ times

2) stay in R_+ $(i \leq L)$ doesn't increase

$\Rightarrow i++$

$\Rightarrow i \leq L$

$\Rightarrow \sum_s (ds + \log n) =$

\Rightarrow occurs $\sum_s ds$ times

$$= \sum_s ds + \sum_s \log n = O(d + m \log n)$$

$$\cdot L = \frac{N}{m} \rightarrow \text{total length of strings} \quad \sigma = |\Sigma|$$

$$\quad \quad \quad \rightarrow \# \text{ of strings}$$

Lower bound

\Rightarrow sorting strings according to 1st character
 $\Omega(n \log_2 n)$ (is sorted)

\Rightarrow ^{Worse +} consider all the distinguish prefix that are in
of all the strings

$$\Rightarrow \Omega(\sum_s d_s) = \Omega(d)$$

$$\Rightarrow \Omega(d + n \log n)$$

MKQ MATCHES THIS BOUND

Let R be set of strings of variable length

$\text{MKQ}(R, i)$

if $|R| \leq 1$

 return R ;

else

p = extract a string to be a "good" PIVOT

$$R_{<} = \{s \in S : s[i] < p[i]\}$$

$$R_{=} = \{s \in S : s[i] = p[i]\}$$

$$R_{>} = \{s \in S : s[i] > p[i]\}$$

$$A = \text{MKQ}(R_{<}, i); B = \text{MKQ}(R_{=}, i+1); C = \text{MKQ}(R_{>}, i)$$

Return $A \cdot B \cdot C$ // (where \cdot is file concatenation)

Cost: Invariant: ~~stays~~ stays in A shores $i-1$ chars

② CASES

\Rightarrow good pivot
mean balanced
partitions

\Rightarrow 1) If stays pos $R_L \vee R_S \Rightarrow$ occurs

$\log_2 n$ ad
① doesn't increase

Stays \downarrow in $R_+ \Rightarrow i+1$

$i \leq L$

occurs $i \leq \sum_s d_s$

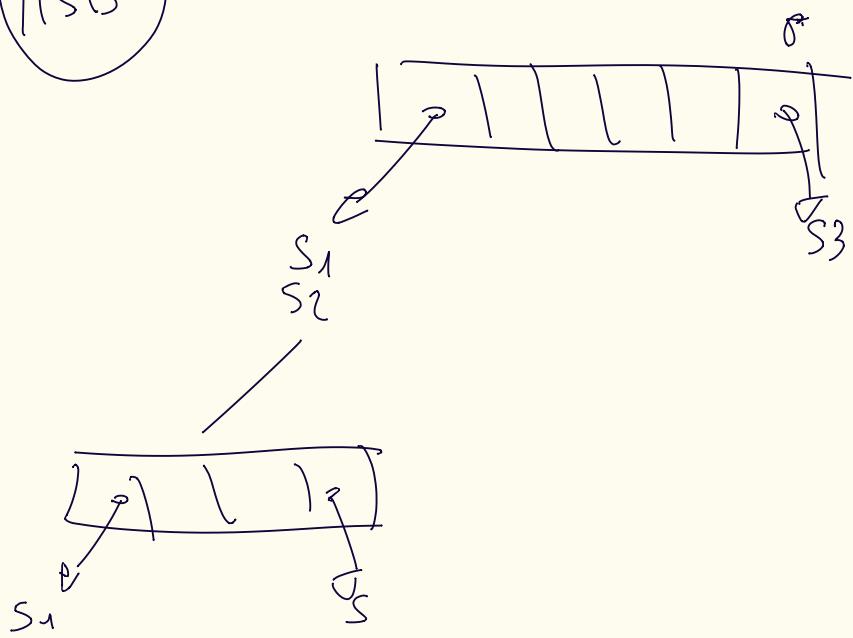
$\Rightarrow O(\sum_s d_s + \log n) = O(d + n \log_2 n)$ ✓

MSD : $O(d \cdot \sigma)$

with hash $O(d + \lg_2 \sigma)$

$$\hookrightarrow \text{SD } O\left(\frac{L}{\lg_2 \sigma}\right)$$

(MSD)



- sort acc. to 1st digit

then
until each
every strip is
in just 1
bucket



Arrays have size σ (at every level)

and I'm "bucketing" the strips (~~strips~~ accordingly
to their distinguish prefixes)

$$\Rightarrow O\left(\sigma \cdot \sum_S d_S\right) = O(\sigma \cdot d) = O(d \cdot \sigma)$$

—
— with hash: every node has a ht \Rightarrow size σ

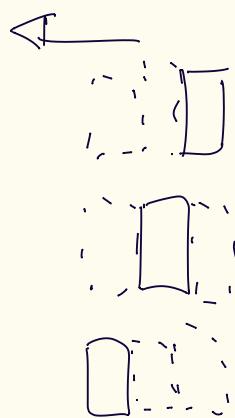
and I've ~~to~~ to take into account all the
edges (they share the algorithm \Rightarrow comparison)

$$O\left(\underbrace{\# \text{edges}}_{\text{en}} \lg \underbrace{\# \text{edges}}_{\text{en}}\right) = O\left(\sum_S d_S \lg \ell_H\right) \leq \\ \leq O(d \cdot \log \sigma)$$

L81)

bucket sort $\mathcal{O}(n + \max \text{rep.})$

as $L=3$



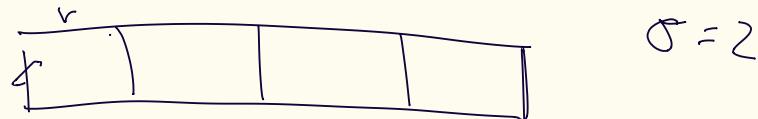
with strip, compare & shift per time

$$\Rightarrow \mathcal{O}(L \cdot (n + \sigma))$$

2 digits

$$\mathcal{O}\left(\frac{L}{2} \cdot (n + \sigma^2)\right)$$

\Rightarrow don't use digits but bits
and compare r-bits at time



$$\sigma = 2$$

$$\Rightarrow \mathcal{O}\left(\frac{b}{r} \cdot (n + \sigma^r)\right) = b = \cancel{2} \text{ bits}$$

$$= \mathcal{O}\left(\frac{b}{r} (n + 2^r)\right)$$

want r that minimize the func

$$2^r \leq n \Rightarrow \text{not the case}$$

$$2^r > n \Rightarrow \frac{b + 2^r}{r} \Rightarrow \frac{b + 2^{\log_2 n}}{\log_2 n} =$$

$$= \frac{\cancel{b + n}}{\log_2 n} = \frac{\cancel{L}}{\log_2 n} \quad \checkmark$$