## String Sorting

$N$ = #strings

$D$ = Total length of distinguished prefix

- variable length = $\ell$ max
- $\tau = |\Sigma|$ alphabet size.

qsort $(A, N, \text{sizeof}(\text{char} *), \text{strcmp})$

↳ thus deploying the power of a cmp-based sort procedure

**PROBLEM**  Time cost = $O(N\ell \log N)$ with many RANDOM mem accesses.
↑ ave-len



long strings, soon diff ⟹ AMORTIZED ANALYSIS
cmp from scratch, always

**NOTE**  $\Sigma \xrightarrow{\text{mapped}} \{0, 1, 2, \dots, \tau-1\}$

strings ⟶ numbers in base $\tau$.

INT. keys
BASIC BLOCK

- KSort $\cong$ Counting sort  (sort only lists)

  1) Keys are integers (SMALL) in $\{0, 1, \dots, K-1\}$  $(K = \tau$ in strings$)$

  2) Key $k$ goes to bucket $b[k]$.

  $\dfrac{\text{Time}}{\text{Space}} = O(n + K)$         OK if $K = O(n)$  ⟹ small alphabet

**Property**  STABILITY  + optimal if →

- Radix Sort $\xrightarrow[\text{iterative}]{}$ LSD ⎡needs all keys have equal length but logical padding
            $\xrightarrow[\text{recurse}]{\text{split}>}$ MSD  ⎣ (strings)

**IDEA**  ① strings as binary sequences   $\underline{0101101 01}$
                                          $\underbrace{4}$
                                          $\overline{16}$   ⟷ base
                                          $256$             TODO →

② strings as sequence of groups of bits, and apply KSort

**LSD**  $d$ = #groups   ⎫ $d(\tau+n)$ Time and $(\tau+n)$ space
        $K$ = max digit symbol  ⎭

$$s: \langle 017, 042, 666, 007, 111, 911, 999 \rangle$$

111 < 017   first round
017 < 111   third round

**proof.** (induction)    $s_{i-1}$ = sorted wrt $(0, ..., i-1)$ digits

We sort by the $i$-th digit, if two numbers have it different
it implies the sorted order and we are done; if it is equal, the
stability preserves the order.

Can we choose $d$ and $K$? We can control them.

$K \Rightarrow$ groups of $\log_2 K$ bits $\Rightarrow d = \dfrac{l}{\log_2 K}$

$$\text{Time} = \frac{l}{\log_2 K}(K + u)$$

→ large $K \Rightarrow$ few groups $\Rightarrow$ costly Ksort
→ small $K \Rightarrow$ many groups $\Rightarrow$ many phases

breakeven point $K = \Theta(u)$, because this is in any case the
cost of one-phase of Ksort.

$$\text{Time} = \frac{ul}{\log_2 u} \qquad \text{Space} = ul$$

Comparison-based
LB = $nl \log n$

Proceeds left→right by distributing keys iteratively

MSD RADIX
counted as a function of D
$O(D + u\sigma)$

Problem 1: Horrend. overheads in the recursion
Problem 2: When buckets are small Ksort is the best

#: small M ⇒ fast !!   since D ≥ u   | σ-partit°

**Multi-key Quicksort**    optimal bound (in the comp-model)

(LB)  $n \log n + D$

**Key idea**  Ternary partition
→ bucket formation is made easy

al | **p** habet
al | **i** gnment
al | **l** ocate
al | **p** orithm
al | **t** ernative
al | **t** ernate
=

→ yet potentially unsorted
← SORTED

Multikey QS $(R, \ell)$      $R =$ set of strings of common prefix

① if $|R| \leq 1$ return $R$      of length $\ell$

② choose pivot $p \in R$

③ $R_< = \{ s \in R \mid s[\ell{+}1] < p[\ell{+}1] \};$

   $R_= = \{ s \in R \mid s[\ell{+}1] = p[\ell{+}1] \}$

   $R_> = \{ s \in R \mid s[\ell{+}1] > p[\ell{+}1] \}$

④ A = Multikey QS $(R_<, \ell)$    (advance)

   B = Multikey QS $(R_=, \ell{+}1)$

   C = Multikey QS $(R_>, \ell)$

⑤ Return $A \cdot B \cdot C$

<br>

- Correctness is obvious
- Cost = Count the number of comparisons To execute Step ③
  - Case $s[\ell{+}1] \neq p[\ell{+}1]$
    - assume (perfect) choice of pivot $\longrightarrow$ R is halved

      (Too strong ⇓ Random)    $\longrightarrow$ Total charge on $s$ is $\log n$

      $\longrightarrow$ # Total comparisons $= n \lg n$

  - Case $s[\ell{+}1] = p[\ell{+}1]$
    - advance of one char $\longrightarrow$ no more than D Total adv

Integer Arithmetic
MSD Radix Sort    What about using $\sigma$ buckets, and not just 3 ?

   deploy INTEGERS

① Distribution Takes $(n_i + \sigma)$ and involves all strings for at most D steps $\left( \sum_{\text{over } n_i} n_i \Rightarrow D \right)$    (OK if $n_i > \sigma$ / Recall for next consideration)

     (not all equal chars)

② Every non-trivial partitioning Takes $O(\sigma)$ Time and this occurs no more than $n$ Times. ( ☛ Non Trivial partition is a Tree with fan-out $\geq 2$ and $n$ leaves).

$$D + n\sigma$$

$$\boxed{\begin{array}{l} \text{Multi key QS} = D + u \log u \\ \text{MSD Radix S} = D + u\sigma \end{array}}$$

Understand algorithms to "combine" them :

① LSD = Ksort + repeat

**now** When the #strings to partition is $\boxed{n_i \leq \sigma}$ use Multikey since it is $\Sigma$-indep.

$$\Rightarrow \boxed{D + u \log \sigma}$$

• We are avoiding the $\sigma$-cost of small-bucket formation which is too much when $n_i \ll \sigma$.

| comp-based | $D + u \log u$ | Multikey QS | OPTIMAL |
|---|---|---|---|
| integer | $D + u \log \sigma$ | MSD + Multikey | |
| constant | $D$  // $D + u\sigma = O(D)$ | MSD Radix | OPTIMAL |

OBSERVATION

Cost Multi-key $\longleftrightarrow$ #distinguishing chars $\Rightarrow$ D
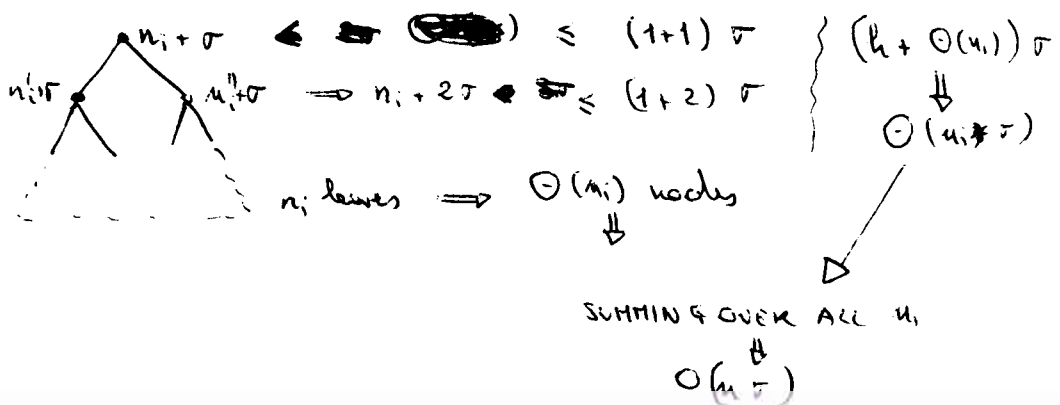
$\hookrightarrow$ #u $\log$ #u

when using COMBINATION. $2°$ term is $\sum u_i \log u_i \leq u \log u$

Cost MSD Radix

• When $n_i \geq \sigma$ $\Rightarrow$ optimal cost $\Theta(u_i) \Rightarrow D$

• When $u_i \leq \sigma$ $\Rightarrow$ poor cost $\Theta(\sigma)$ which we pay at each exit.

$$n_i + \sigma \leq \qquad ) \leq (1+1)\sigma \qquad \Big\{ (h + O(u_i))\sigma$$
$$n'_i \sigma \qquad u''_i + \sigma \Rightarrow n_i + 2\sigma \leq (1+2)\sigma \qquad \Updownarrow$$
$$\Theta(u_i \neq \sigma)$$

$n_i$ leaves $\Rightarrow \Theta(n_i)$ nodes
$\Downarrow$

SUMMING OVER ALL $u_i$
$\Updownarrow$
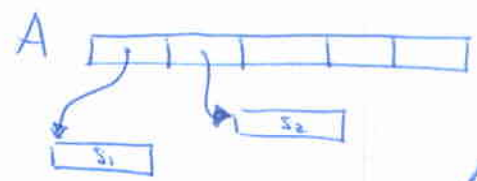$O(u \cdot \sigma)$

$\boxed{\text{Lecture \#3}}$  Strings in a set

- Sorting → Search data structures
  ↳ Duplicate removal
  ↳ Join in DBs.

  OTHER NOTES

- Search + Storage $D = \{ abaco, abba, paolo, pattern, patty \}$

  ① Array of string pointers, randomly allocated strings

  A
  

  ○ $P \times \log N$ Time cost / $\lg N$ I/Os / Space : $L_{TOT} + N + N$
    ↳ random memory accesses

  ○ Interpolation Search ⟹ $\log \log N$ → $mid = \dfrac{low + (K - A[low])}{A[high] - A[low]}$

    $\left(\text{Assumes random distribution of keys in the range } [A[low], A[high]), \text{ given that there are } (high-low) \text{ keys}\right)$

    $\cdot (high - low)$

    - String ≡ number in base $\sigma$
    - recent $\Delta$ = gap ratio = $\dfrac{\max (x_i - x_{i-1})}{\min (x_i - x_{i-1})}$

    ① bin $B_i$ represents a range of size $\dfrac{x_u - x_1}{u}$ (equally part')

    ② binary search on $B_J$ where $J = \dfrac{K - x_1}{x_u - x_1}$

    Cost is $O(\log \Delta)$ → $\Delta$ = polylog $u$ per unit dist.

    Clearly $\max x_i - x_{i-1} \geq \dfrac{x_u - x_1}{u}$ (avg distance)

    $\min x_i - x_{i-1} \overset{def}{=} \dfrac{\max x_i - x_{i-1}}{\Delta} \geq \dfrac{x_u - x_1}{u \Delta} \Rightarrow |B_i| \leq \Delta$

- How To save $\approx \log_2 B/\ell$ cache/block misses
  and be output-sensitive?
  → write strings contiguously
  → ~~░░░~~ [~~░~~ ↳ last $\log_2 B/\ell$ accesses are in the
  same block, both for pointers & strings.

- How To save space?
  → ~~●~~ - ~~●~~   ~~░░░░░░░~~
  ⟹ Blocking : reduce #pointers $N \longrightarrow \dfrac{L}{B}$  [1 byte more or 10]
  $\overset{\text{CONS}}{\longmapsto}$ This requires Block scanning but page already fetched
  $\overset{\text{PRO}}{\longmapsto}$ Front-coding    (35÷40%)
  $$\text{Time} = \log_2 \frac{L}{B} + 1 \quad I/Os \quad \boxed{\ell \le B}$$
  $$\text{Space} = \underset{\underset{\text{strings}}{\uparrow}}{L_{TOT}} + \underset{\underset{\text{pointers}}{\nwarrow}}{\frac{L}{B} 4} ~~\blacksquare~~ + \underset{\nwarrow 10}{N}$$

- Can we speed-up prefix-searches?  [TRIE]
  → uncompressed : one node per char
  → compressed : squeeze out the unary path. [need SKIP or DIGIT NUM]



  ① P $I/Os$  (indep. N)
  ② ~~░░░~~  ~~░░░~~
     12÷20 N bytes

  case strings start all with same digits

#internal nodes ≤ N-1 , #leaves = N

- branching impacts.
  ① binary impl = leftmost child + sibl ⟹ 2 pointers × internal node
     but $O(\sigma)$ branch time
     ~~BE BACK ON THIS~~ (circled)
  ② hash table ⟹ $O(1)$ avg branch time, $O(N)$ space
  ③ array ⟹ $O(1)$ branch time, $O((N-1)\cdot\sigma)$ space !!

Time/IO saving

## What about I/O-issues? | Two-level memory

1^level
- ~~Node~~ implementation impacts on the I/O-performance
  → larger ~~blocks~~ → fewer cached → more I/O-misses.

- Given 2-levels → Trie on a set sample
                  ↘ scan on the rest (buckets)



① TRADE-OFF: ↑ Smaller buckets → bad FC but faster search

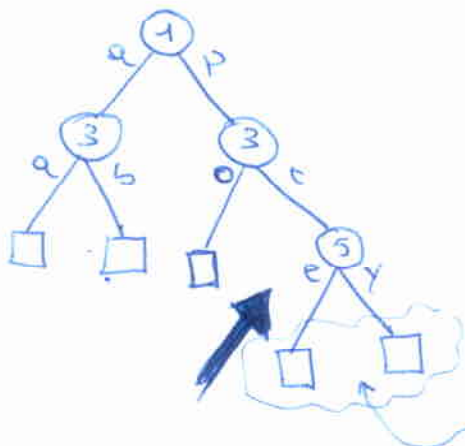② COPIED SPACE: Much because we need to keep entire strings to solve edge labels

■ Two improvements
   - Trie → drop edge labels + strings → more "copied" strings
                                        ↘ difficult search
   - Buckets + FC → no bucketing but FC on all (more robust)

## BLIND TRIE | described in the attached paper but it is not described how to lex-search [MY NOTES]
                                                        ↳ crucial to find where to proceed.



P = patricia
   ↑ ↑↑
   1 3 5

1^phase (either we reach a leaf or pick one descending)

lcp(P, "patty") = pat|E| ⟹ P < patty

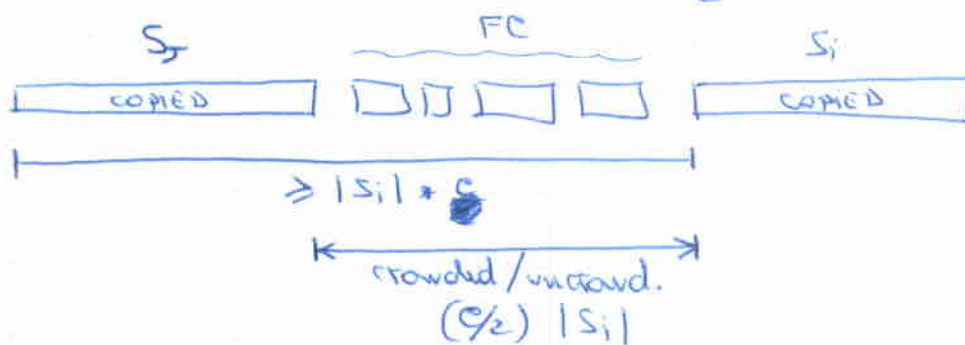## SPACE | proportional to sampled number and not to length
                                    (thousands of strings in L2)

$\boxed{\text{TIME/IO}}$   1 string check $\Big\langle$ buckets (obvious)
                                    $\swarrow$ LPFC (cool idea!!)

① Front-code a string $S$ iff the cost of decoding it $\leq c|S|$

[# previous explicit char ~~...~~ ~~...~~ to scan in order to decode]

② Decoding procedure is not affected

$\boxed{\text{TEO}}$ $\forall e > 2$, space $= (1+\varepsilon)$ FC, Time $\dfrac{|S|}{\varepsilon}$, where $\varepsilon = \dfrac{2}{e-2}$

$S_j$                    FC                    $S_i$

[ COPIED ] [ ⬚⬚ ⬚ ⬚ ] [ COPIED ]

$\geq |S_i| * \dfrac{c}{2}$

crowded / uncrowd.
$(c/2) \, |S_i|$

•) Decompose the sequence of copied into $(\text{UNCROWD}) \cdot (\text{CROWD})^*$

•) if $S$ uncrowded, iT is preceded by at least $\dfrac{c}{2}|S|$ chars of FC

•) if $S_i$ crowded, then $|S_j| \geq \dfrac{c}{2}|S_i| \implies |S_i| \leq |S_j| \dfrac{c}{2}$

(geometric)

RUN LEN
~~...~~ $< |\text{UNCROWD}| \cdot \displaystyle\sum_{x=0}^{\infty} \left(\dfrac{c}{2}\right)^x = |\text{UNCROWD}| \cdot \dfrac{1}{1 - \frac{2}{e}}$

•) Charge this cost onto the $\dfrac{c}{2}|\text{UNCROWD}|$ chars that precede the run

$\implies \dfrac{2}{e-2}$ per char cost $\implies$ this is $\varepsilon$

(first run has no preceding chars but uncrowded is $\textcircled{FC}$-chars)

$\boxed{\text{NOTE}}$ This allows also to drop completely the bucketing
and this result more robust.