

CSE4502/5717
BIG DATA ANALYTICS
EXAM2 REVIEW

SYLLABUS:

- ① Out of Core Graph algorithms
- ② Randomized alg. for selection on a single BKK
- ③ Parallel BKKs Model
- ④ Suffix trees - Applications
- ⑤ Suffix Arrays

① Graph Algorithms: $G(V, E)$.

DFS in a SOURCE DISK.

Graph is supplied as adjacency lists.



Let $d(u)$ be the degree of node u .

$$\begin{aligned} \text{I/O Complexity} &= \sum_{u \in V} \left\lceil \frac{d(u)}{B} \right\rceil \\ &\approx \sum_{u \in V} \left(\frac{d(u)}{B} + 1 \right) = O\left(\frac{|E|}{B} + |V| \right). \end{aligned}$$

② Rand. Alg. for selection.

Input: $X = k_1, k_2, \dots, k_n, (1 \leq i \leq n)$.

Output: The i^{th} smallest of X .

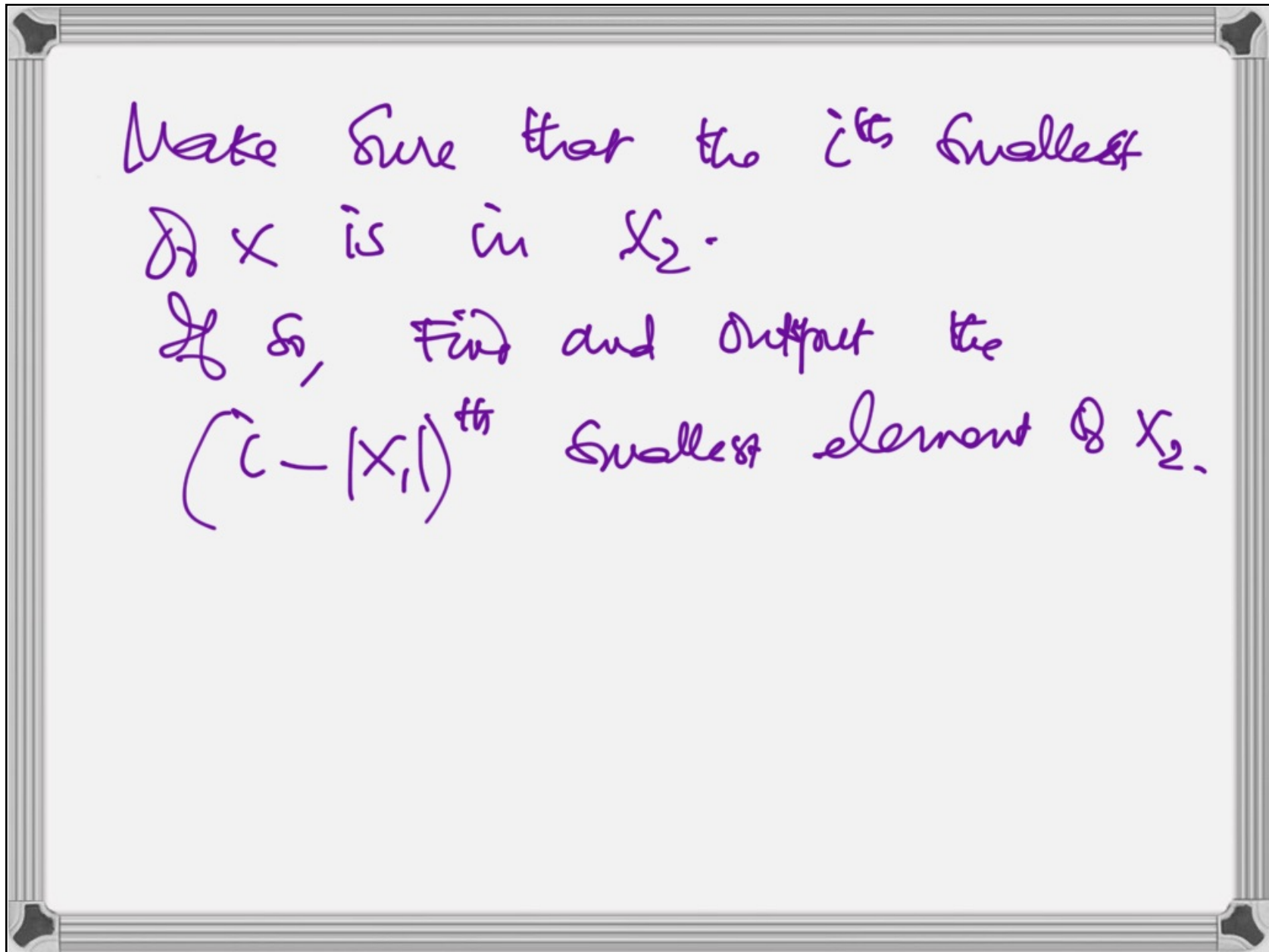
IDEA- Pick a Random Sample S
from X . let $|S| = 8$.

Identify 2 elements in S s.t.

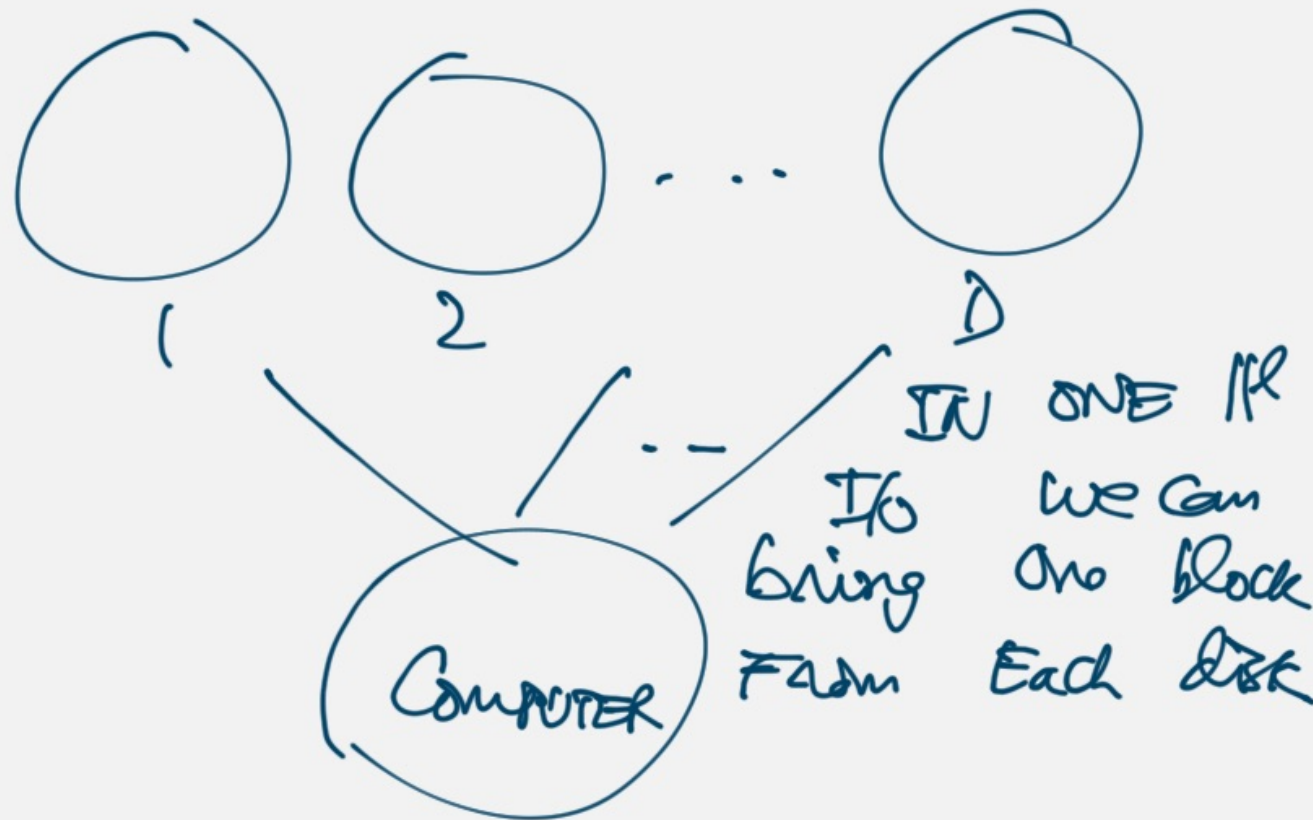
(1) The i^{th} smallest $\exists x \in [l_1, l_2]$
 where l_1 and l_2 are the 2 elements

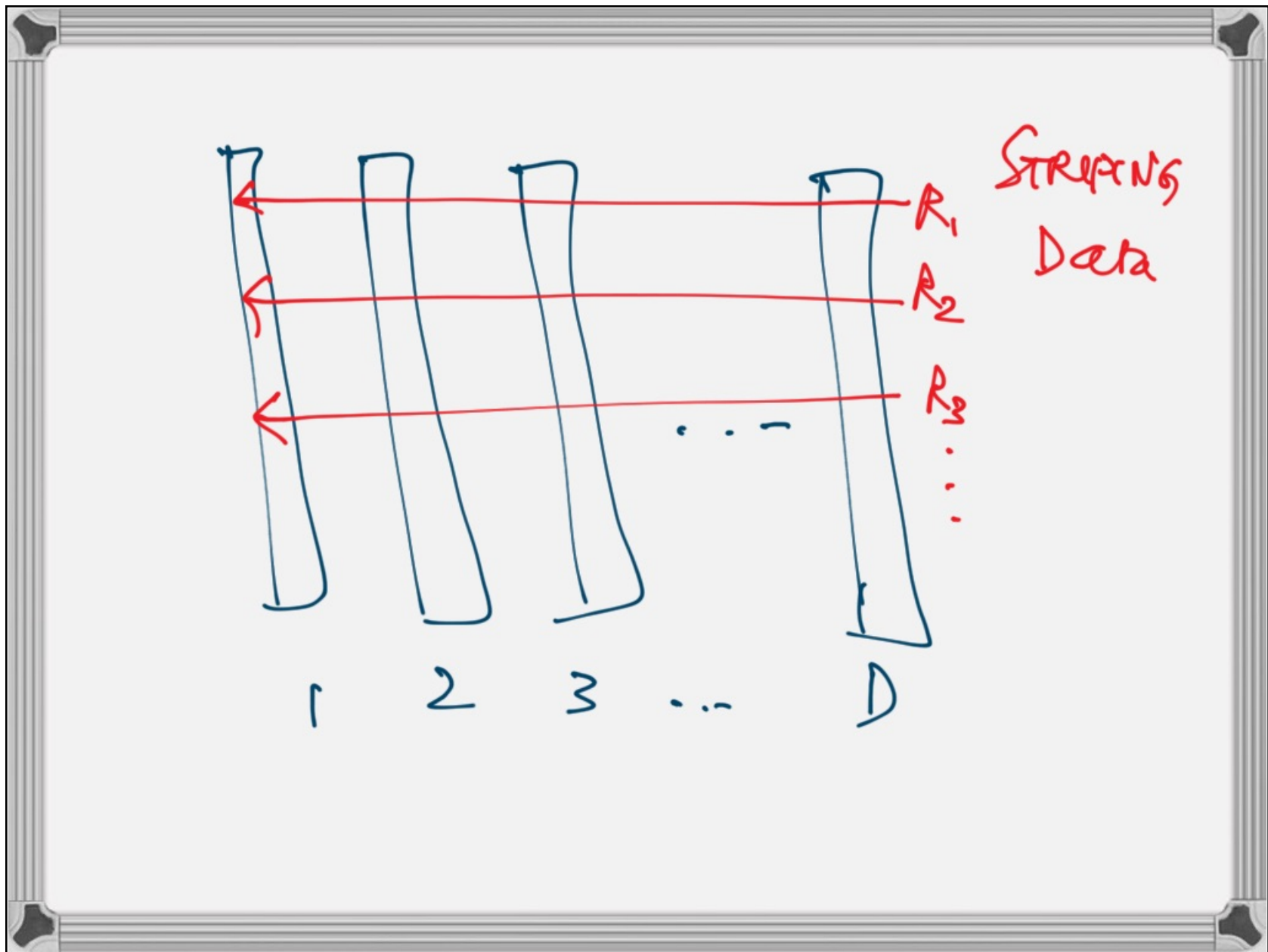
and (2) $\left\{ \{e \in X : l_1 \leq e \leq l_2\} \right\}$ is "small".

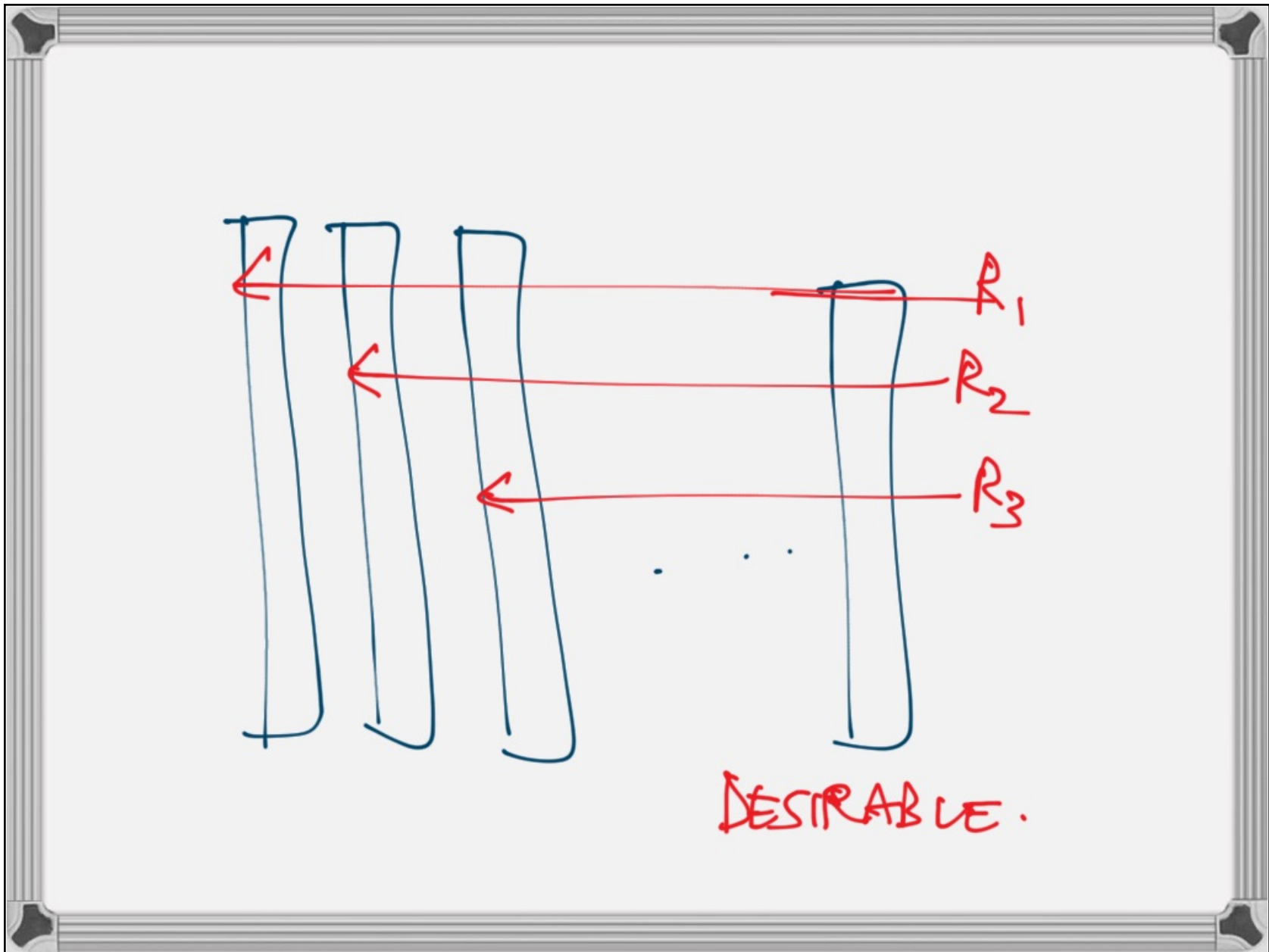
let $X_1 = \{e \in X : e < l_1\}$
 $X_2 = \{e \in X : l_1 \leq e \leq l_2\}$.



Parallel Disks Model:







① Form runs of length M each.

② Merge these $\frac{N}{M}$ runs.

Using k -way MERGE for
some suitable k .

DISK STRIPED MERGE SORT:

Treat the D disks as a single
disk with a block size of BD .

The whiteboard contains handwritten notes in green and red ink. On the left, there are two diagrams illustrating the merge process. The first diagram shows two vertical bars, each with a curved top, representing two sorted sub-arrays. The second diagram shows a single, taller vertical bar with a curved top, representing the merged array. To the right of these diagrams, the text reads:

Use $\frac{M}{BD}$ -way Merge.
 IR I/O Complexity
 $= O\left(\frac{\log(N/M)}{\log\left(\frac{M}{BD}\right)}\right)$ REVS.

Below this, the text reads:

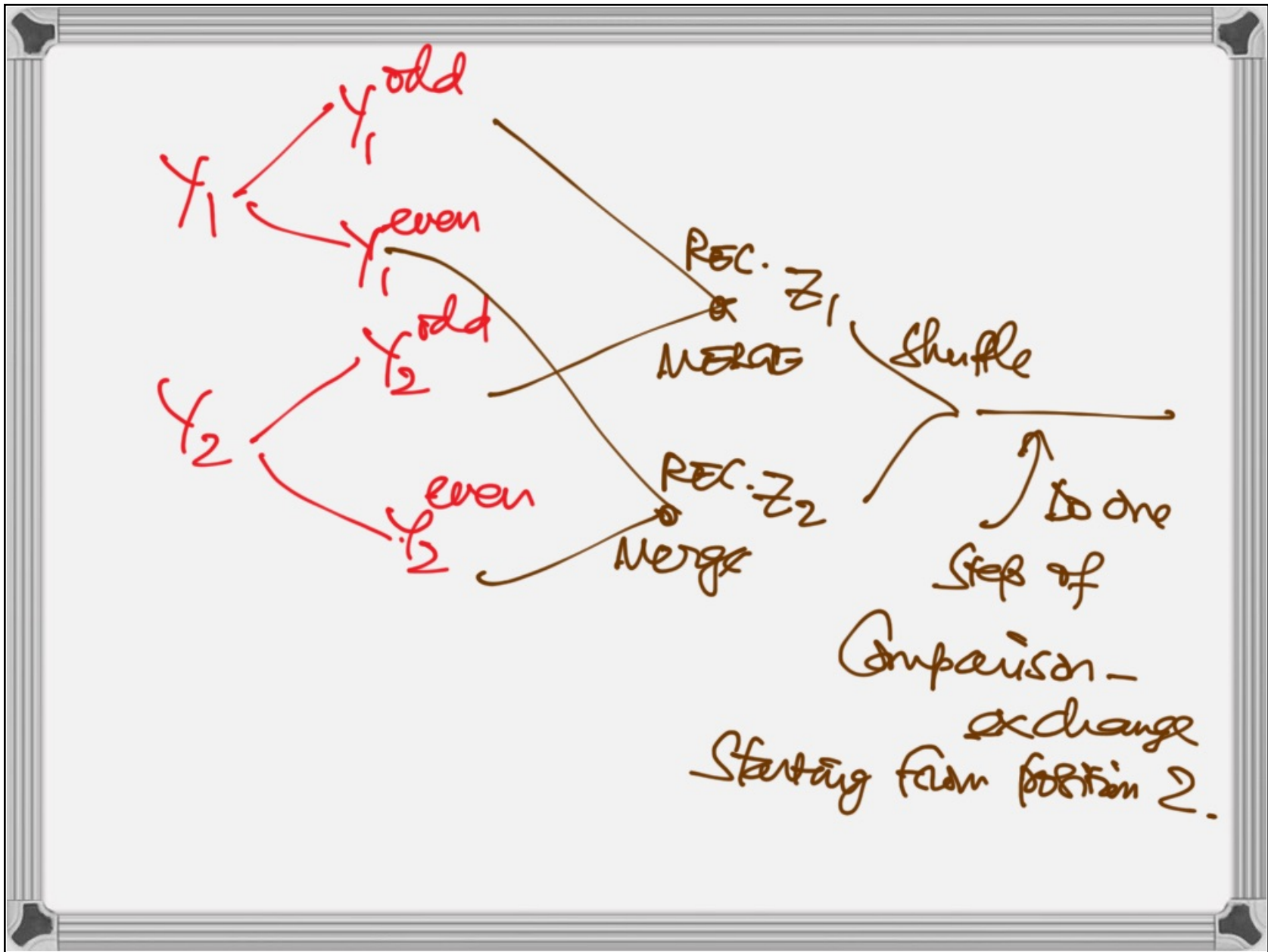
Lower Bound for Sorting: $O\left(\frac{\log(N/M)}{\log(M/B)}\right)$ ←

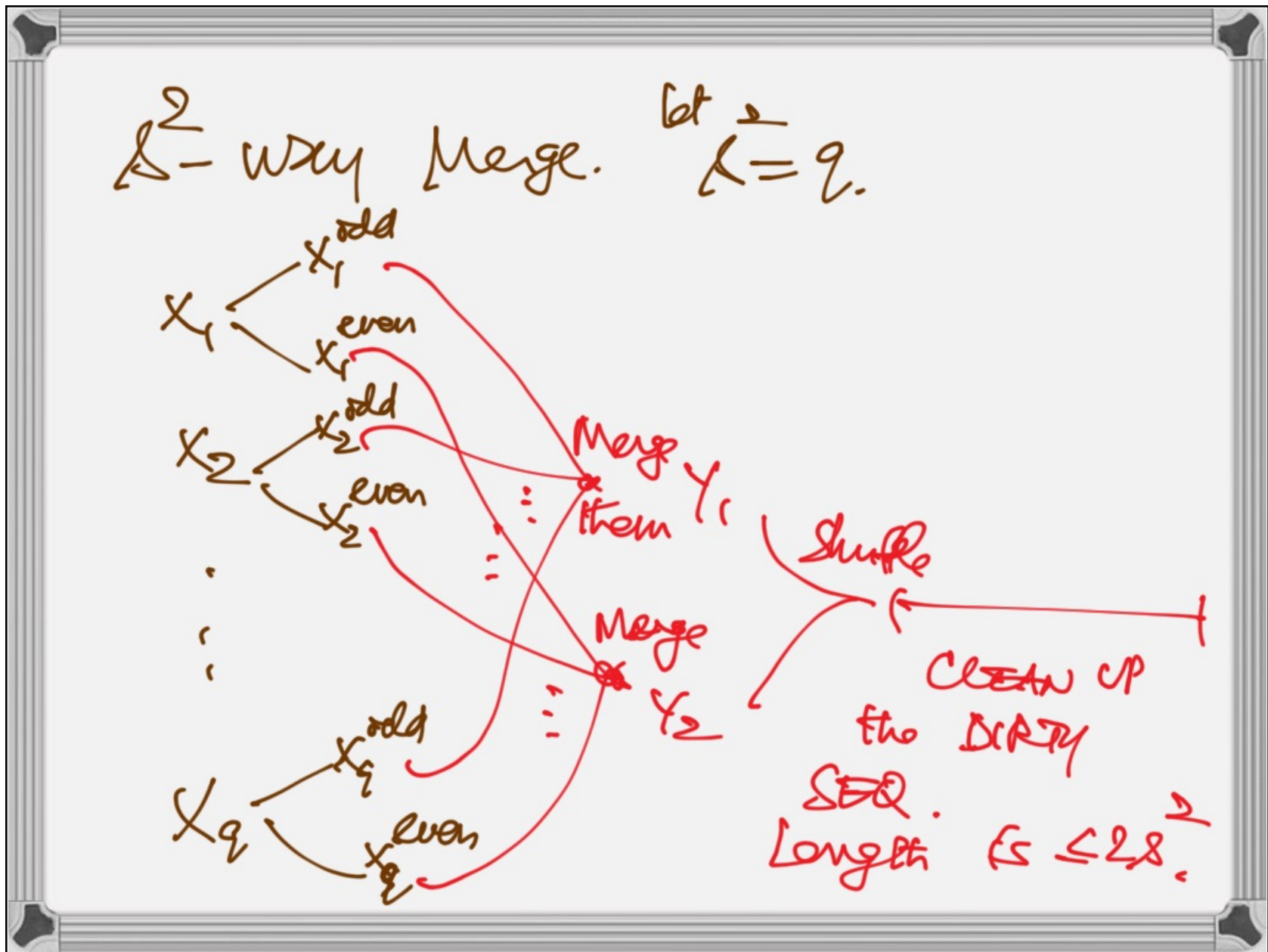
odd-even Merge sort

To sort $X = k_1, k_2, \dots, k_{\frac{n}{2}}, k_{\frac{n}{2}+1}, \dots, k_n$

(1) Sort X_1 recursively; sort X_2 recursively

(2) Merge the sorted runs using the odd-even merge alg.



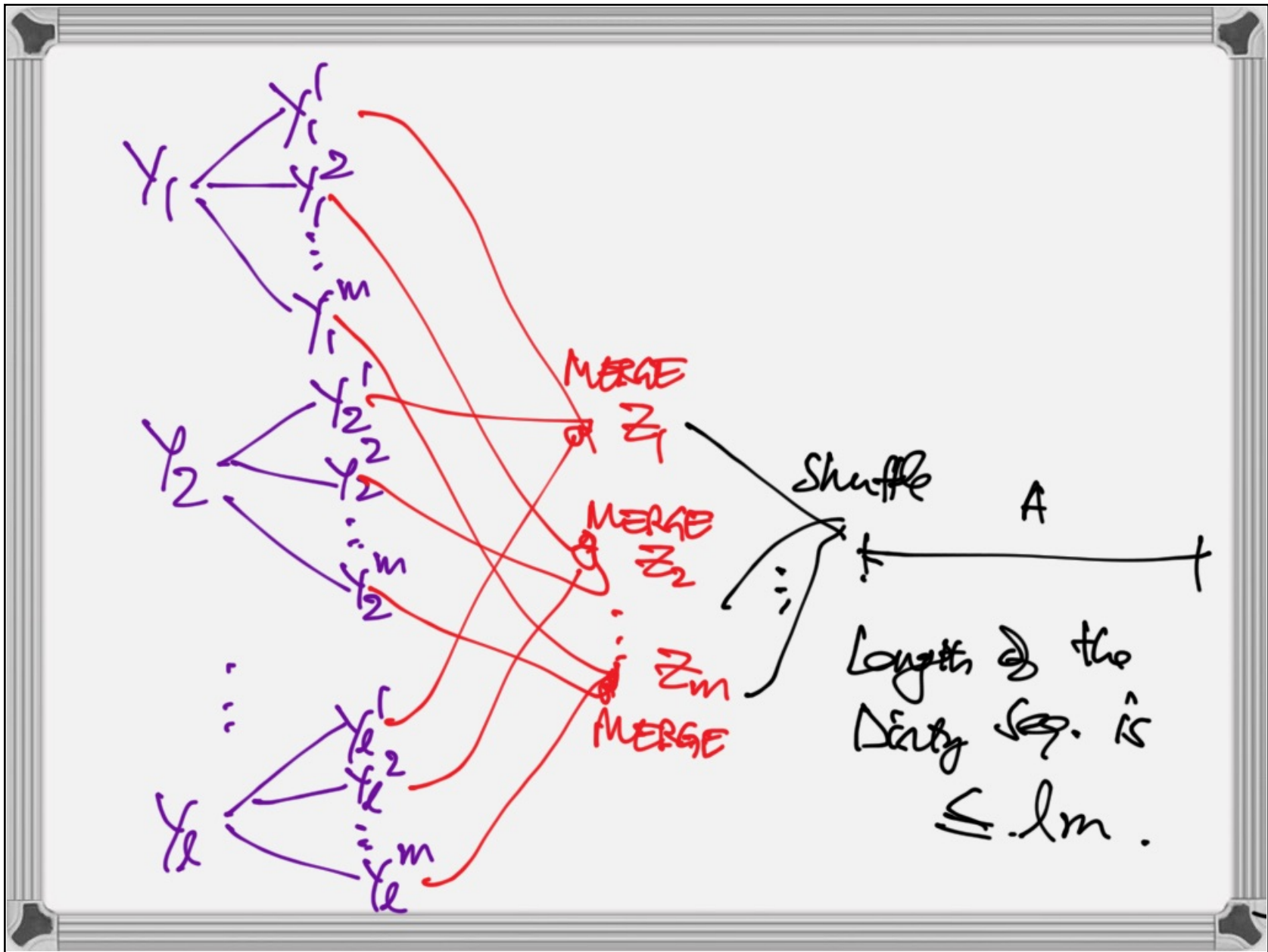


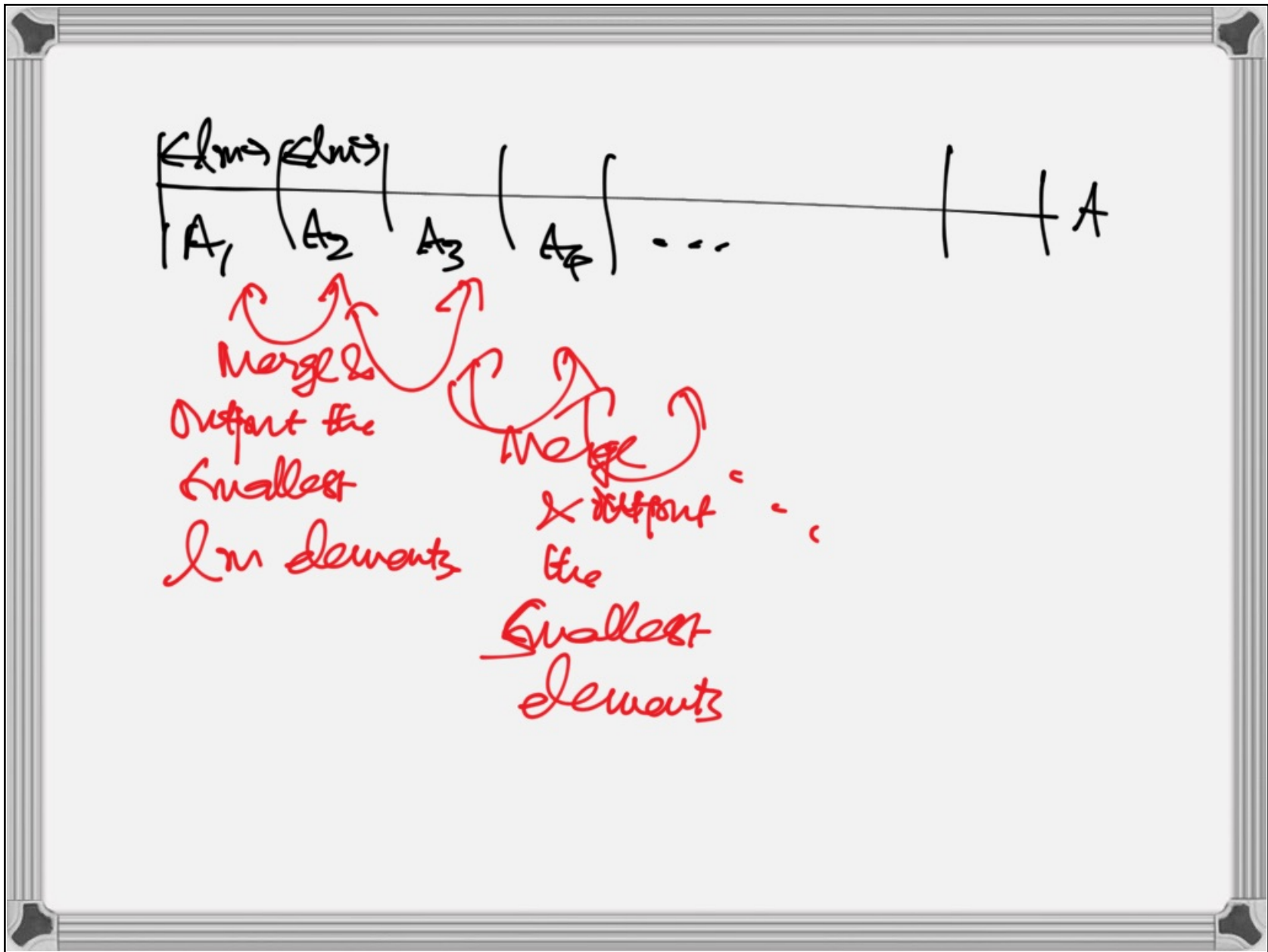
(d, m) -MERGE SORT.

INPUT: $X = K_1, K_2, \dots, K_n$

Output: Sorted X .

- ① Partition X into X_1, X_2, \dots, X_l
 $|X_i| = \frac{n}{l}, \quad 1 \leq i \leq l.$
- ② Sort EACH X_i recursively, $1 \leq i \leq l$
 Let these sorted sequences be $Y_1, Y_2, \dots, Y_l.$
- ③ Merge Y_1, Y_2, \dots, Y_l using the (d, m) -MERGE algorithm.





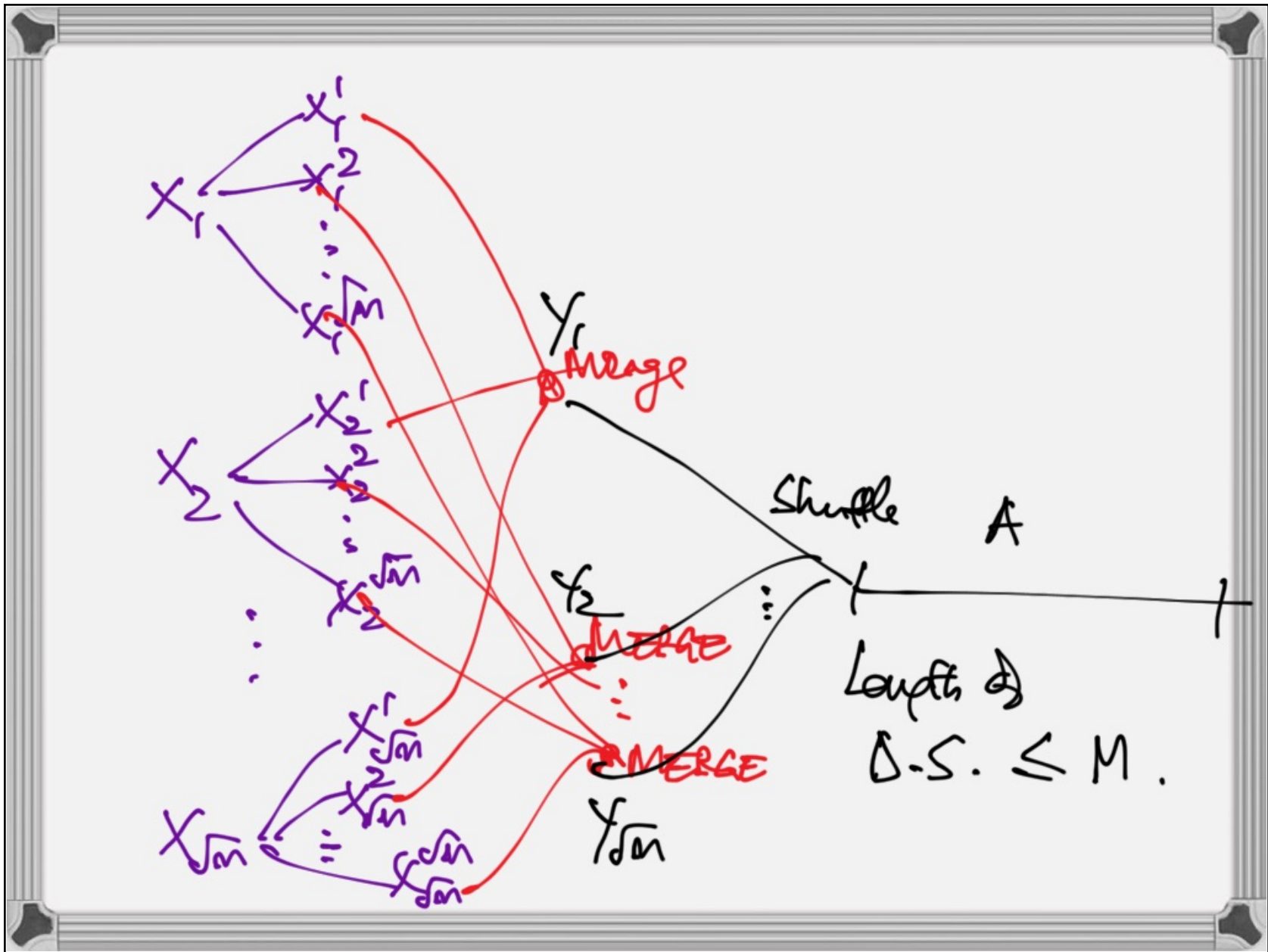
Parallel I/O Complexity

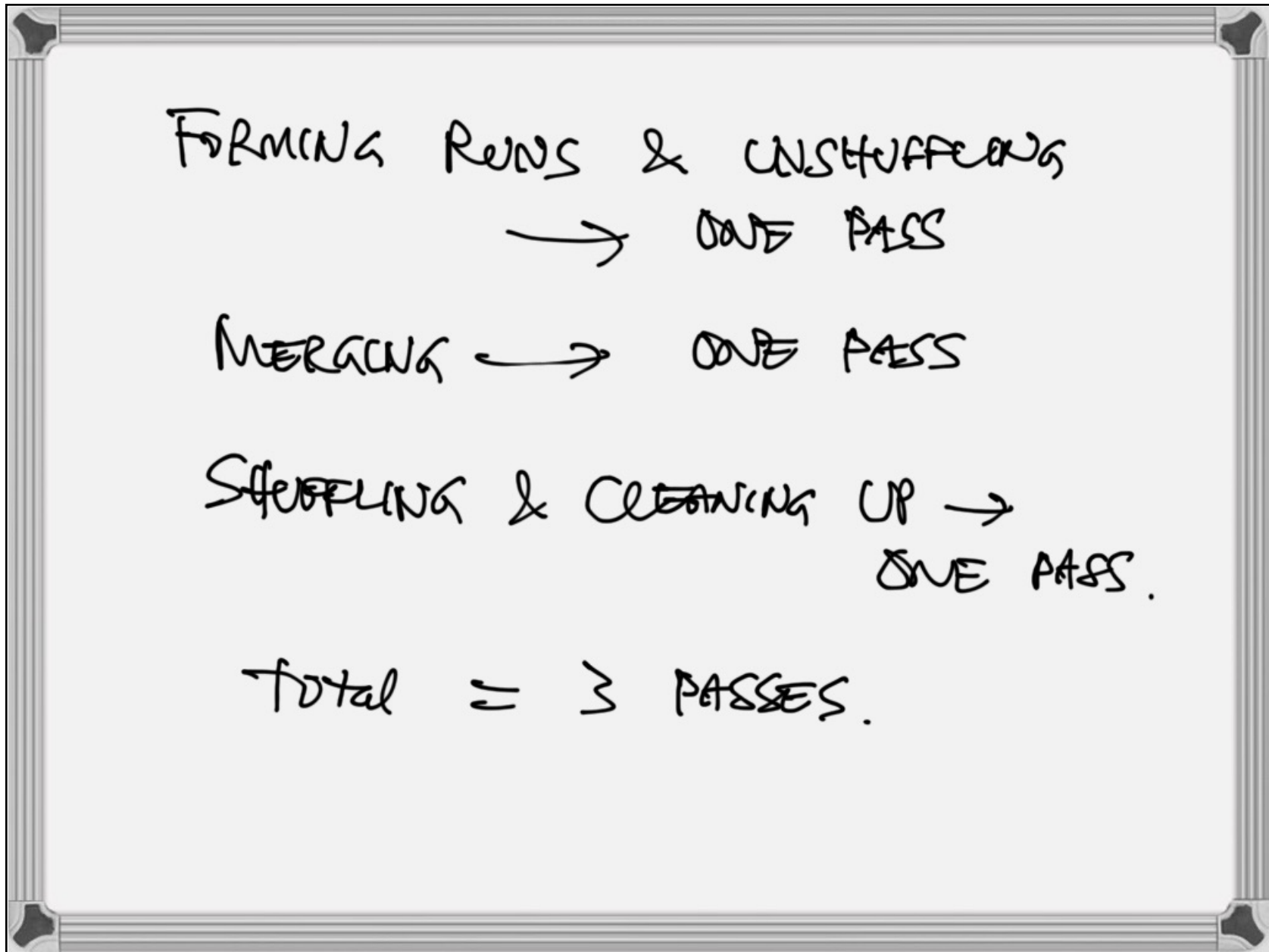
$$= \left[\frac{\log(n/m)}{\log(\min\{\sqrt{m}, \frac{m}{b}\})} + 1 \right]^2$$

Example: Sort $M \times M$ elements
where $B = D = \sqrt{M}$.

① Form Rows of length M
EACH. let the runs be
 $X_1, X_2, \dots, X_{\sqrt{M}}$

② Merge these runs using
(l, m)-MERGE alg with
 $l = m = \sqrt{M}$.





SUFFIX TREES.

$$T = a_1 a_2 \dots a_m \in \Sigma^m$$

Suffix tree on T is a rooted tree S.T.

- ① \exists a leaf for every suffix of T
- ② Every internal node has a degree ≥ 2

- ③ Every edge is labelled with a substring of T
- ④ For any node the labels on no two outgoing edges can start with the same character
- ⑤ We can label the leaves with integers in the range $[1, m]$.
- ⑥ The concatenation of the edge labels starting from the root

and ending in a leaf labelled \bar{c}_i , will be the i^{th} suffix of T
 i.e. $a_1 a_2 \dots a_m$.

We can construct a suffix tree on T in $O(m)$ time.

We always put a $\$$ @ the end of every string, where $\$ \notin \Sigma$.

PROBLEMS:

① STRING MATCHING.

INPUT: $T = t_1 t_2 \dots t_m$ $m \gg n.$
 $P = p_1 p_2 \dots p_n$

Output: All the occurrences of
 P in T .

Construct a suffix tree Q on T .
Start from the root & try to
find a match for P along a
path.

Run Time = $O\left(\begin{matrix} m + k \\ + n \end{matrix}\right)$ \swarrow # of matches.

Given multiple strings S_1, S_2, \dots, S_q

We can construct a G.S.T. \mathcal{Q}

on S_1, S_2, \dots, S_q in time

$$O\left(\sum_{i=1}^q |S_i|\right).$$

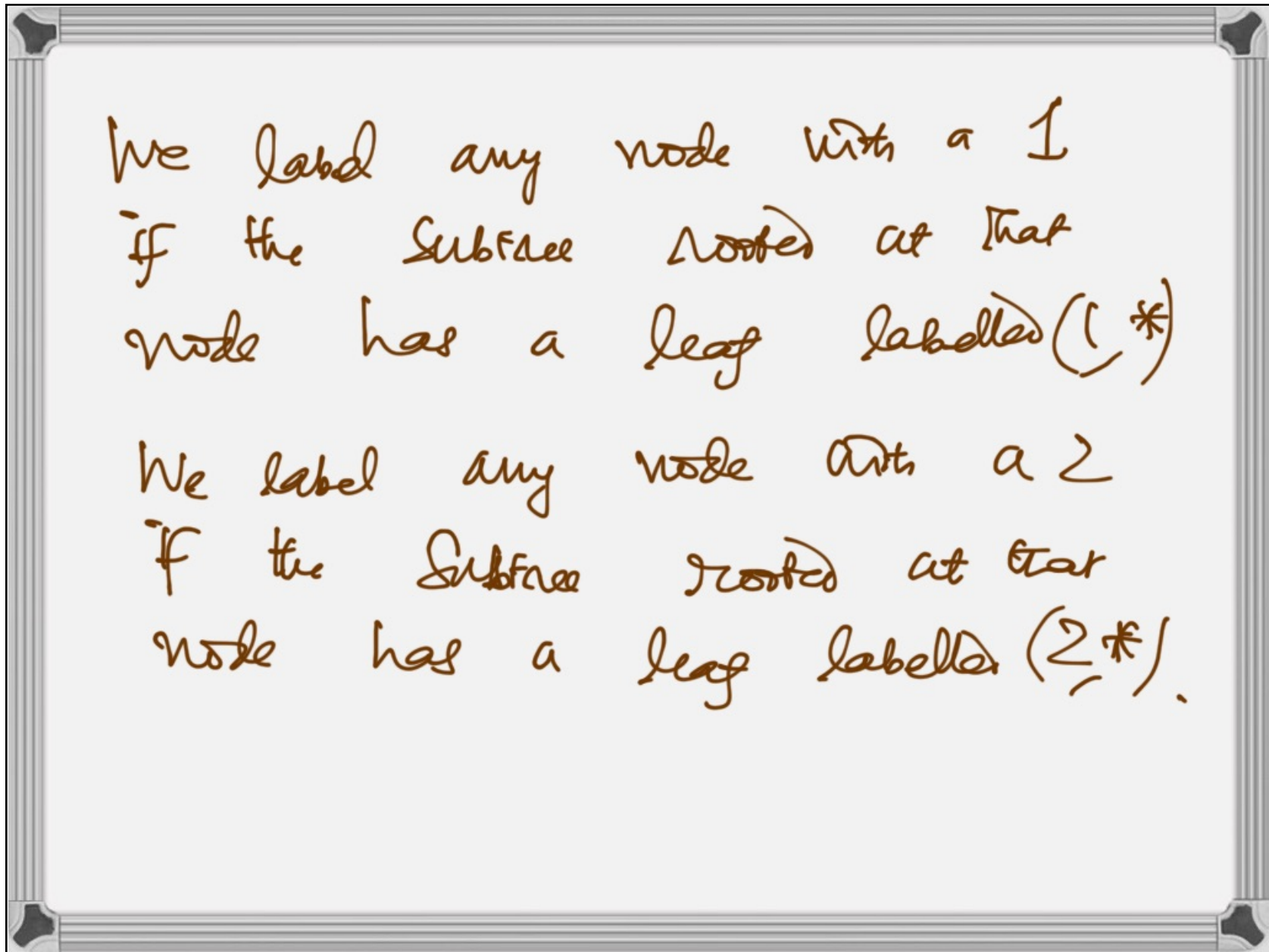
Every is labelled with a pair (i, j) .
 \Rightarrow This leaf corr. to the j^{th} suffix
 \mathcal{D} S_i .

LONGEST COMMON SUBSTRING:

INPUT: S_1, S_2 .

Output: LCS of S_1 & S_2

IDEA. Construct a G.S.T. Q
on S_1 and S_2 . \Rightarrow TIME = $O(n^2)$
 $\text{if } |S_1| = |S_2| = n.$



Traverse \mathcal{Q} one more time to look for the node u that is labelled with 1 and 2 and whose facing depth is the largest. Output the path label of u .

Real time = $O(n)$.

All PAIRS PREFIX-SUFFIX PROBLEMS

INPUT: S_1, S_2, \dots, S_n

Output → The longest prefix of S_j
that is a suffix of S_i
For every i and j .

FACT. We can solve this in
 $O(Mn^2)$ time, where $M = \sum_{i=1}^n |S_i|$.



$\tilde{c} \in L[u]$ if
 u has a terminal
 edge corr. to S_i .

Suffix Arrays:

$T = t_1 t_2 \dots t_m$

Suffix Array for T is an
 Array $SA [1:m]$ of integers,
 s.t. $SA [i]$ is the starting
 position^{in T} of the i th smallest suffix
 of T .

STRING MATCHING USING A SUFFIX
ARRAY:

P

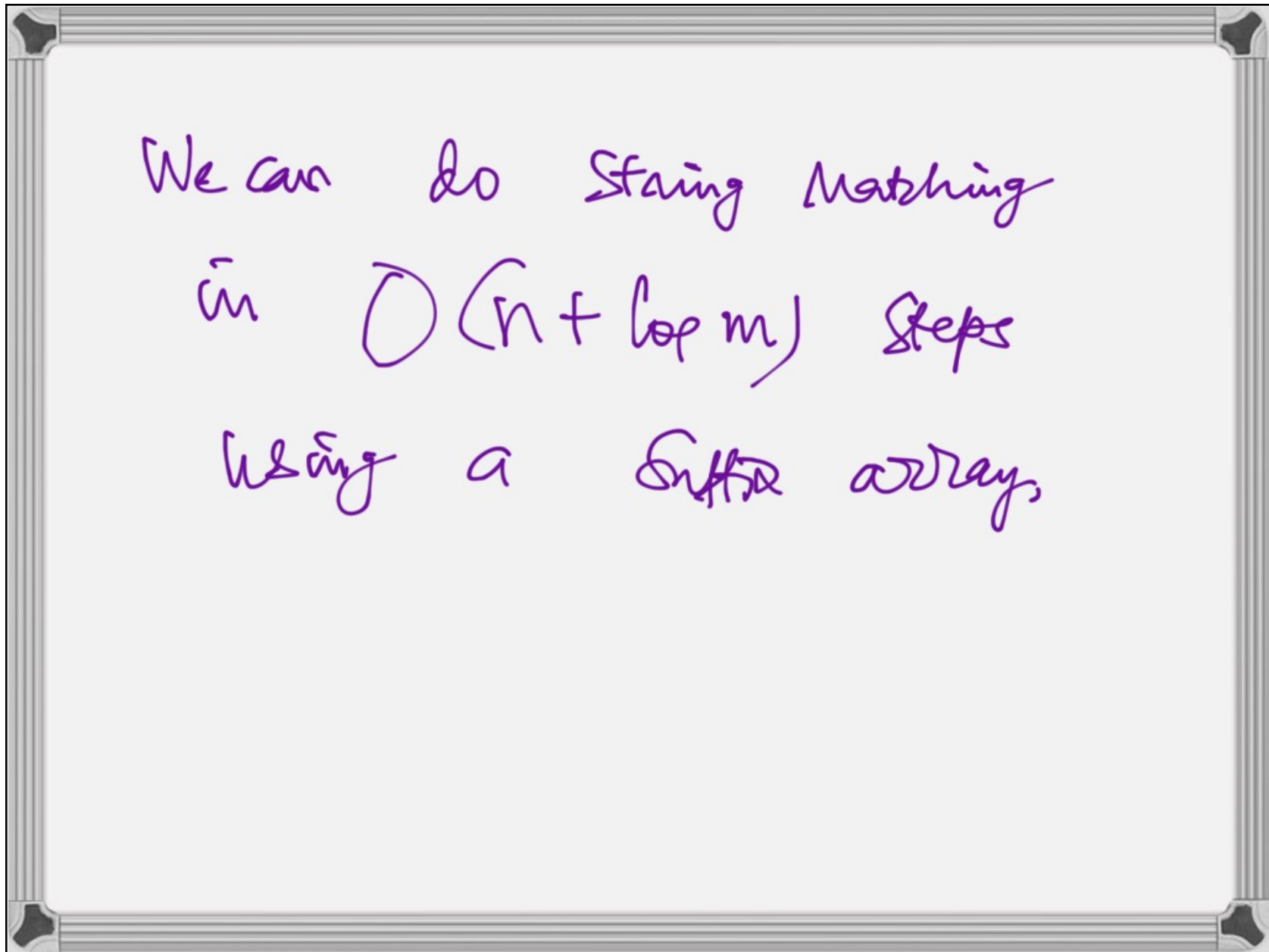
TIME FOR
BINARY
SEARCH
 $= O(n \log m)$

Example: $\Sigma = \{a, b, c, d\}$

$T = cabdacb$

arr

| | | |
|---|---|---------|
| 1 | 2 | abdacb |
| 2 | 5 | acb |
| 3 | 7 | b |
| 4 | 3 | bdacb |
| 5 | 1 | cabdacb |
| 6 | 6 | cb |
| 7 | 4 | dacb |



We showed that a suffix array
on T can be constructed in
 $O(m)$ TIME.

Model Exam 2s

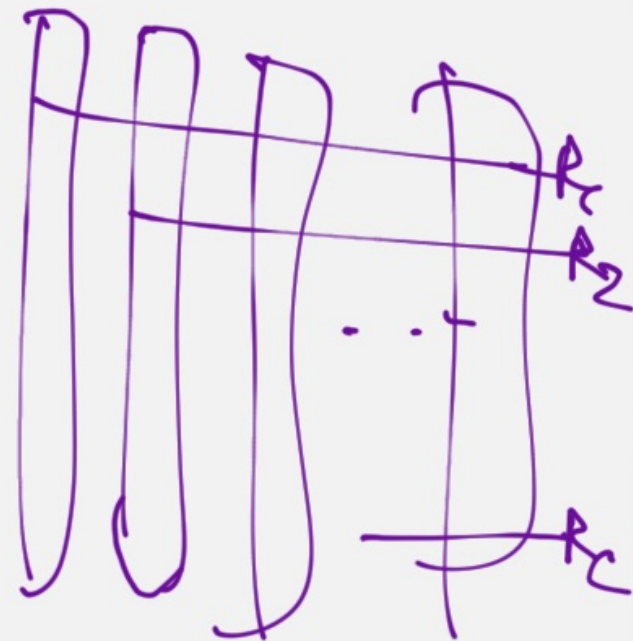
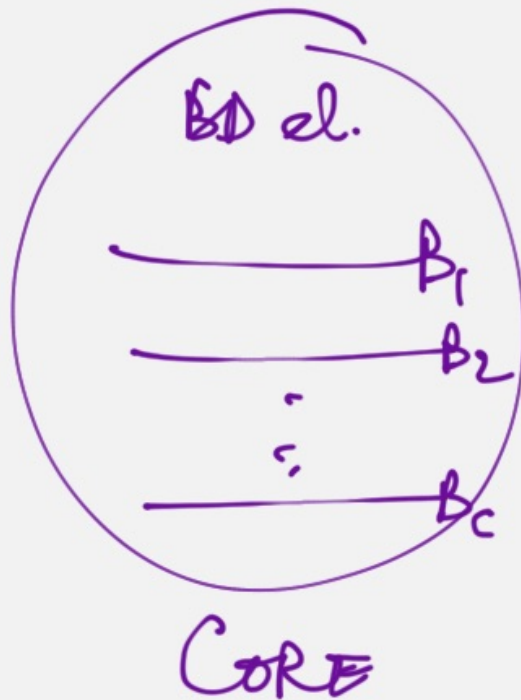
INPUT: $X = K_1, K_2, \dots, K_n.$

of distinct elements = $C.$

Sort in two passes.

IDEA: Bring B elements @
a time.

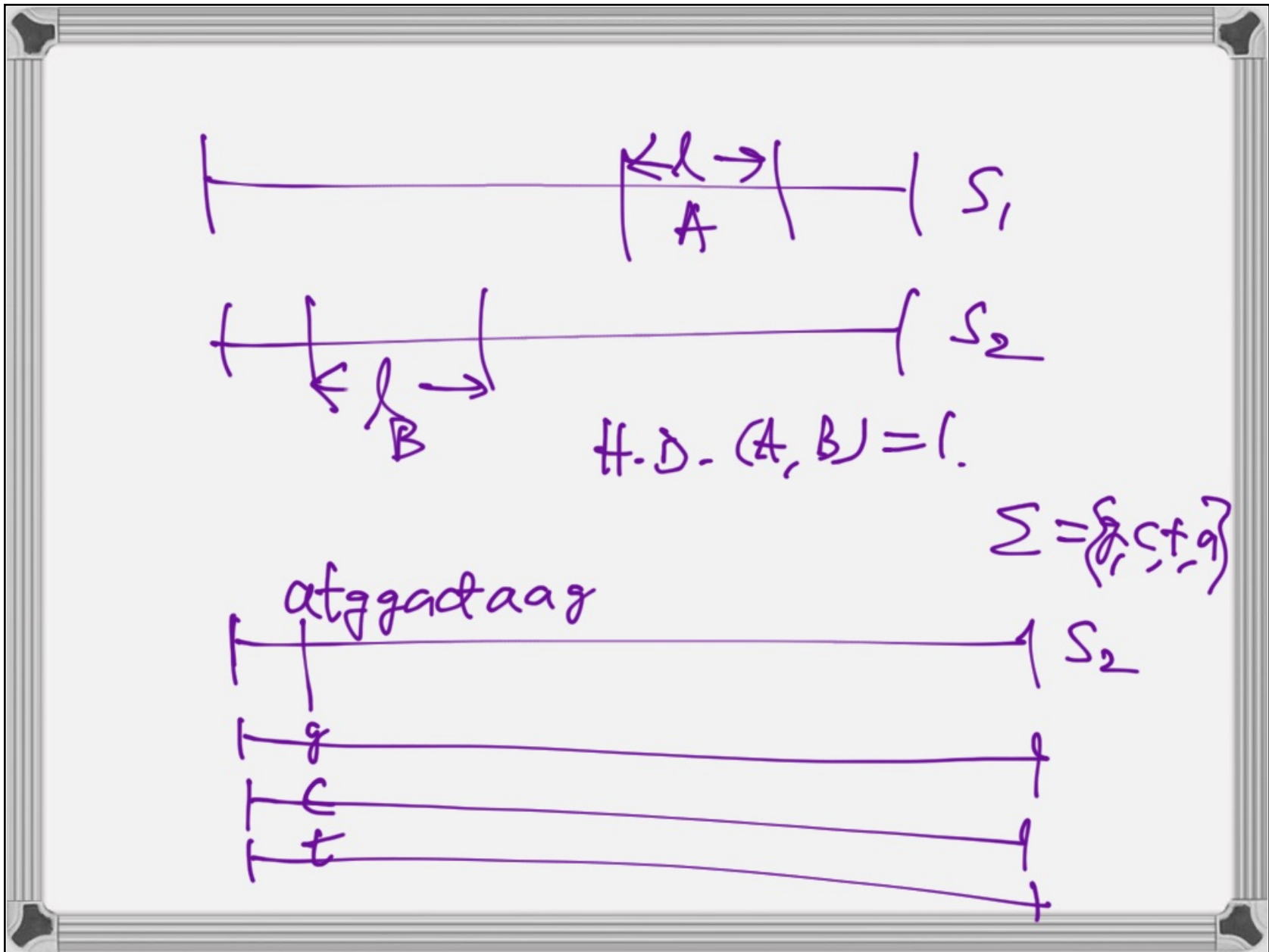
GROW C Runs one for each
distinct value.

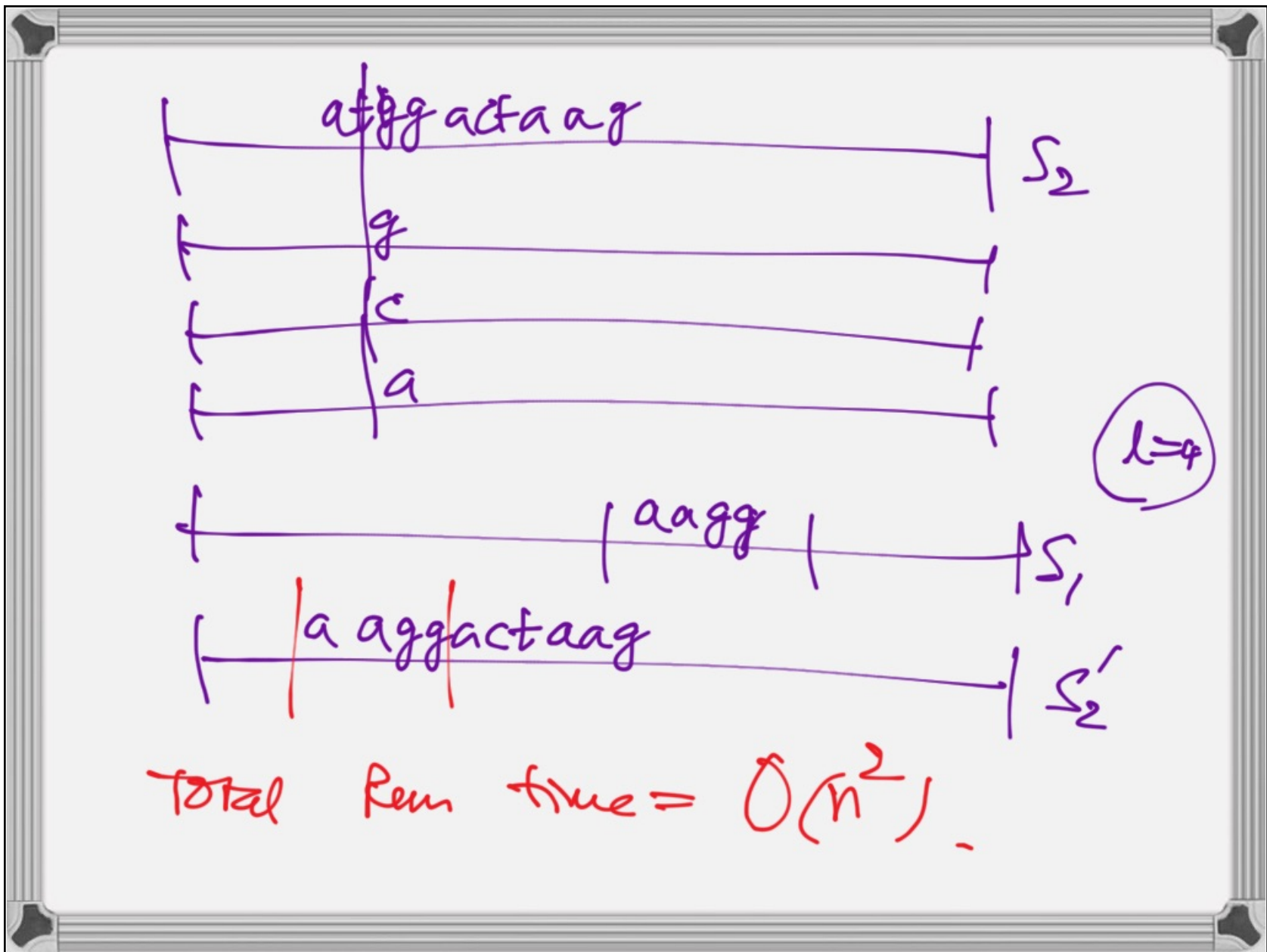


PROBLEM 2: Selection.

In $O(n)$ passes Find the next
Smallest BD elements &
delete them from the input.

Do the above again & again
until the BD smallest we have
contains the i^{th} smallest
Do a selection on these elements





PROBLEM 4.

CONSTRUCT A F.S.T. ON S_1, S_2, \dots, S_k .

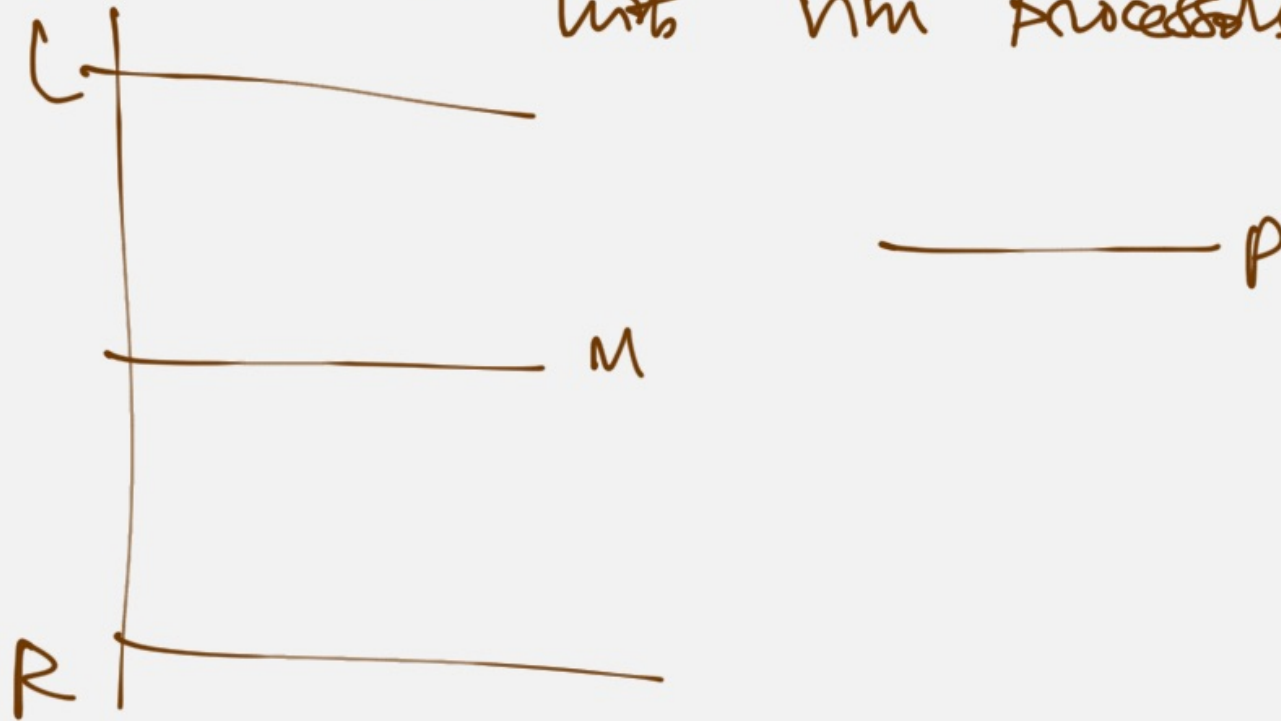
TIME = $O(M)$.

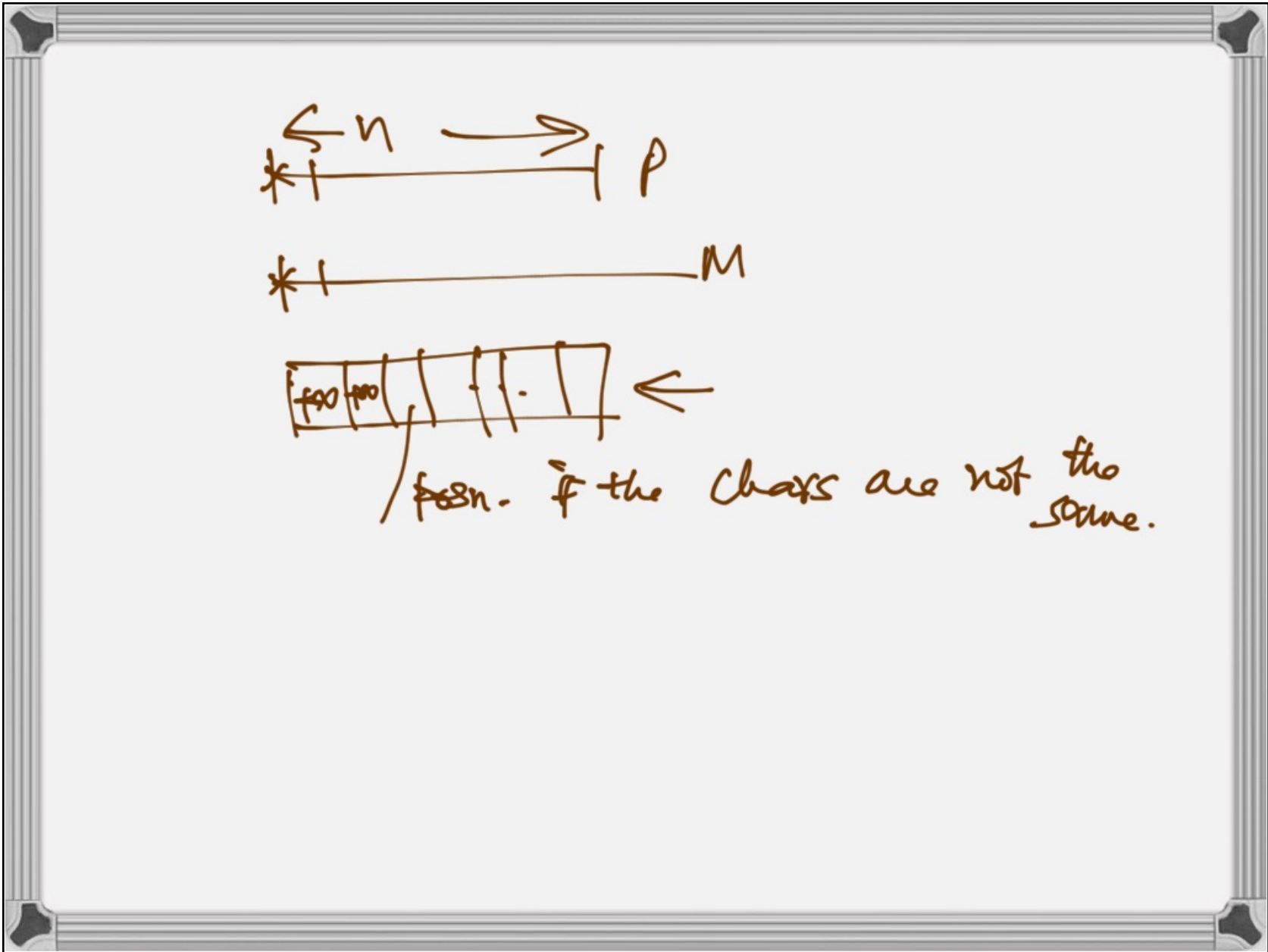
For $1 \leq i \leq k$ do

Traverse Q . Label a node N with i if its subtree has a leaf covr. to S_i

Traverse Q to find the node
 u that is labelled $1, 2, \dots, k$
& whose string depth is
the largest. Output the
path label of u .
Total time = $O(kM)$.

PROBLEM 5. (5a) Scheduling Mapping
into VM processes.





(56) $\rho = n\sqrt{m}$.

The diagram illustrates a process flow with a vertical line intersecting a series of horizontal lines. On the left, a vertical bracket indicates a segment of length \sqrt{m} . This is followed by another bracket of length \sqrt{m} , then three vertical dots, and finally a last bracket of length \sqrt{m} . Each horizontal line segment is labeled in red as "n proc". A red bracket on the right side groups the three horizontal lines corresponding to the three vertical dots.

THE END OF REVIEW!