

ADB Lecture 3

• Introduction to Query Processing

→ Consider a Select query with Search Key $\langle \text{age}, \text{sal} \rangle$

* The way we process the query will depend on the underlying File Organization.

• Heap File

→ Can be an unclustered index (B+ tree or hash)

• Sorted File

→ Can be a clustered index (B+ Tree)

* The query may involve one of many operations

Scan

EqualitySearch

RangeSelection

Insert

Delete

• Fetch all records

• Field = ...

• Fetch the block then locate record

• Field $>, \dots$ (or $\leq, \geq, <, >$)

• Fetch the page • do the operation then write it

→ The key to good Physical design is accurately describing the expected workload

→ workload is a mix of queries and updates
(Search query) (action query)

→ For each query in the workload, need to identify

accessed relations

attributes retrieved TC

attributes in select/join conditions

+ their selectivity
• output length
table length ↓

→ For each update

type of update

attr. affected

may need to update the index as well depending on that.

• Choice of Indexes

- ① → which relations should have indexes? • Freq. of choose
- ② → what fields should be the search key?
- ③ → How many indexes for the relation?
- ④ → what's the type of each?
 - Clustered? • At most 1
 - Hash/Tree?

• A Possible Approach

- Given how DBMS evaluates queries and creates query evaluation plans.
 1. Consider most important queries
 2. Consider the best evaluation plans for them given the current indexes
 3. If a better plan is possible with an additional index, then create it.

Should as well consider the impact on updates before creating an index.



⇒ Index Selection Guidelines

1. For index keys, consider attributes in WHERE

• Hash Index ← Exact Match Condition

• Tree Index ← Range Queries

B+
• Many changes

ISAM

* In general, Clustered Index is helpful for range queries and exact match if there are duplicates

Composite

2. Use multi-attribute Search Keys if WHERE Clause Contains Several Conditions

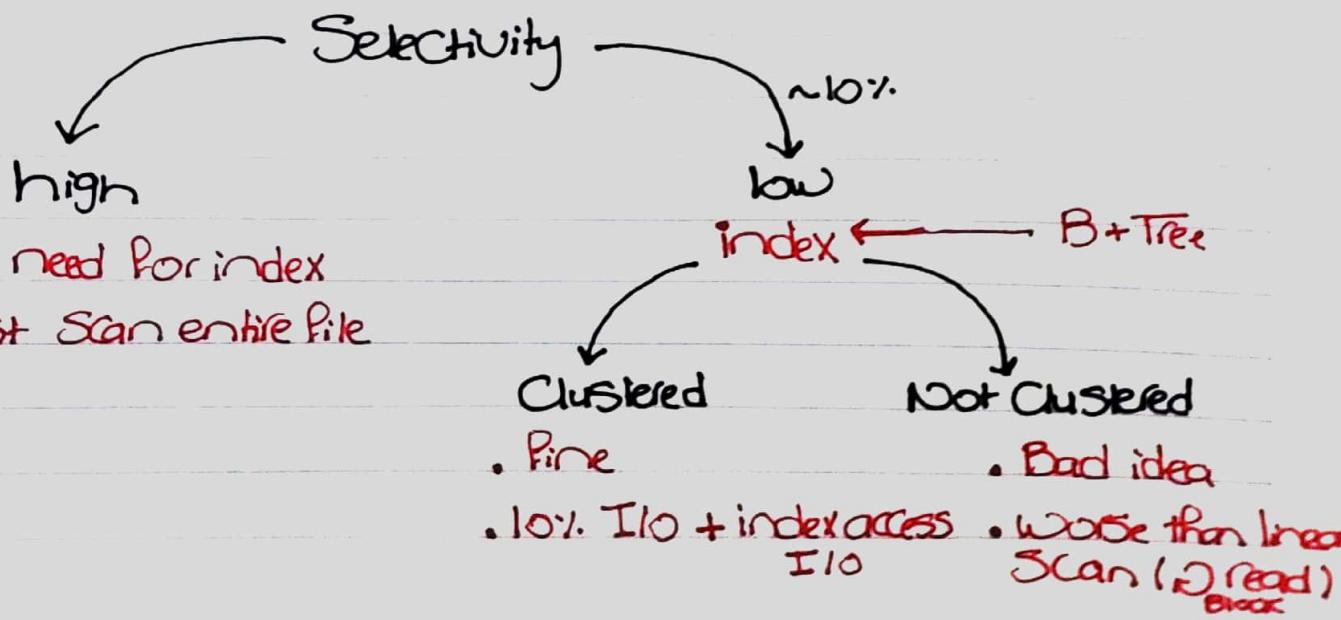
- In Case Range Queries are expected, Order the attributes based on them.

3. Choose indexes that benefit as many queries

- E.g., we know we can have at most one clustered index per relation
→ Choose it based on frequent queries that would benefit the most from clustering.
- Should carefully decide when it comes to clustered index (Expensive updates; although cheaper than sorted files)

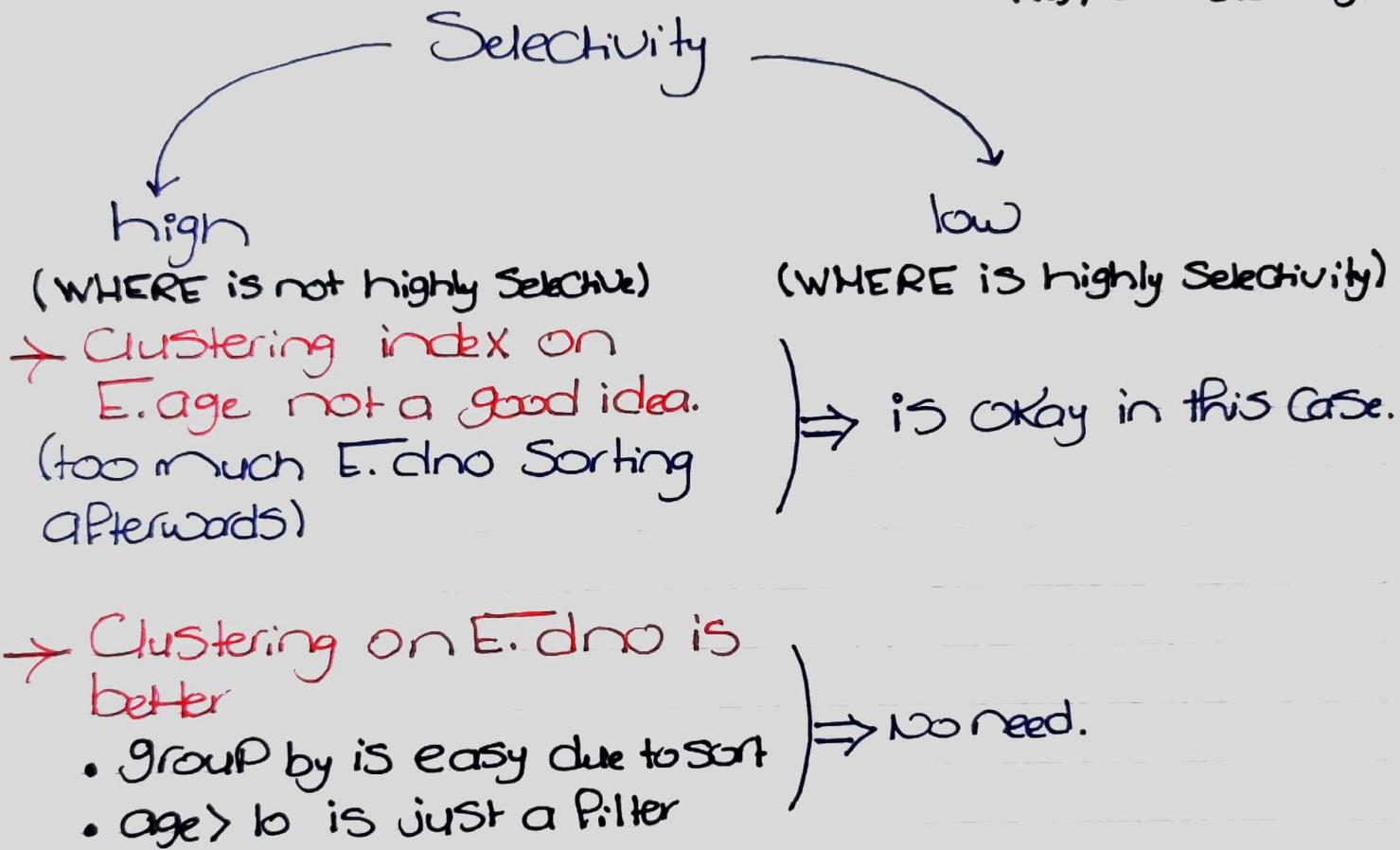
• Examples

① SELECT E.dno
FROM EMP E
WHERE E.age > 40 # Range Query



② SELECT E.dno, Count(*)
 From EMP E
 WHERE E.age > 10
 Group by E.dno

, Suppose an index
 is to be used.
 (e.g., also equality)



③
 SELECT E.dno
 From EMP E
 Where E.hobby = 'Stamps'

} • A secondary index on
 E.hobby can help
 (duplicates are possible)

④
 Select E.dno
 From EMP E
 Where E.eid = 50

} • eid is unique
 • By Guideline 1, hash is
 better this time.

5

Select *
From EMP E
Where Sal > 10

Visually

age	Sal
...	

- all occurrences
in file sorted by age
then Sal.

Consider a Composite index

$\langle \text{age}, \text{Sal} \rangle$ $\langle \text{age} \rangle$

- Not helpful
(sorted by age
then by Sal)

$\langle \text{Sal}, \text{age} \rangle$ $\langle \text{Sal} \rangle$

- Both work
• Sal must be
Prefix

Consider the following modifications

age = 30 and Sal = 4000

→ $\langle \text{age}, \text{Sal} \rangle$

- better than $\langle \text{age} \rangle$ or
 $\langle \text{Sal} \rangle$ (can binary
search twice)
- better than $\langle \text{Sal}, \text{age} \rangle$
if age = ... is more
selective. ①

20 $\langle \text{age} \leq 30 \text{ and } \text{Sal} \leq 4000 \rangle$

→ $\langle \text{age}, \text{Sal} \rangle$ or $\langle \text{Sal}, \text{age} \rangle$

- , intuitive as range queries
benefit much from sorting

age = 30
and
300 $\langle \text{Sal} \leq 4000 \rangle$

→ $\langle \text{age}, \text{Sal} \rangle$

- ① definitely
holds

* Index Only Plans

→ Some queries can be answered just using the
index (no tuples retrieved from relation).

→ in this case, clustering is not important
(connection to relation)

Examples

Select E.dno, Count(*)
From Emp E
Group by E.dno

- Key is E.dno
 - e.g., Count Pointers out of each E.dno value
- Secondary index (non-key)

Select E.dno, MIN(E.Sal)
From Emp E
Group by E.dno

- Key is $\langle E.dno, E.Sal \rangle$
- For each group of dnos, Salary Sorted by default.

Select avg (E.Sal)
From Emp E
Where E.age = 25 and
E.Sal between 3000 and 5000

- Key is $\langle E.age, E.Sal \rangle$
- Find $\langle 25, x \rangle$ then
Keep reading.
3000 or more if not.

* External Sort

- Needed to create clustered indexes (sort w.r.t key)
- 1st step in bulk loading B+tree (instead of repeated insertions, can sort first)
- Eliminating duplicates
- Needed in Sort-merge join algorithm

Problem to Solve

→ Sort 1GB of data with 1MB of RAM while minimizing disk access.

Two-way Sort

→ uses 3 buffers (2 input, one output) • size = block
→ Suppose file has 2^k pages

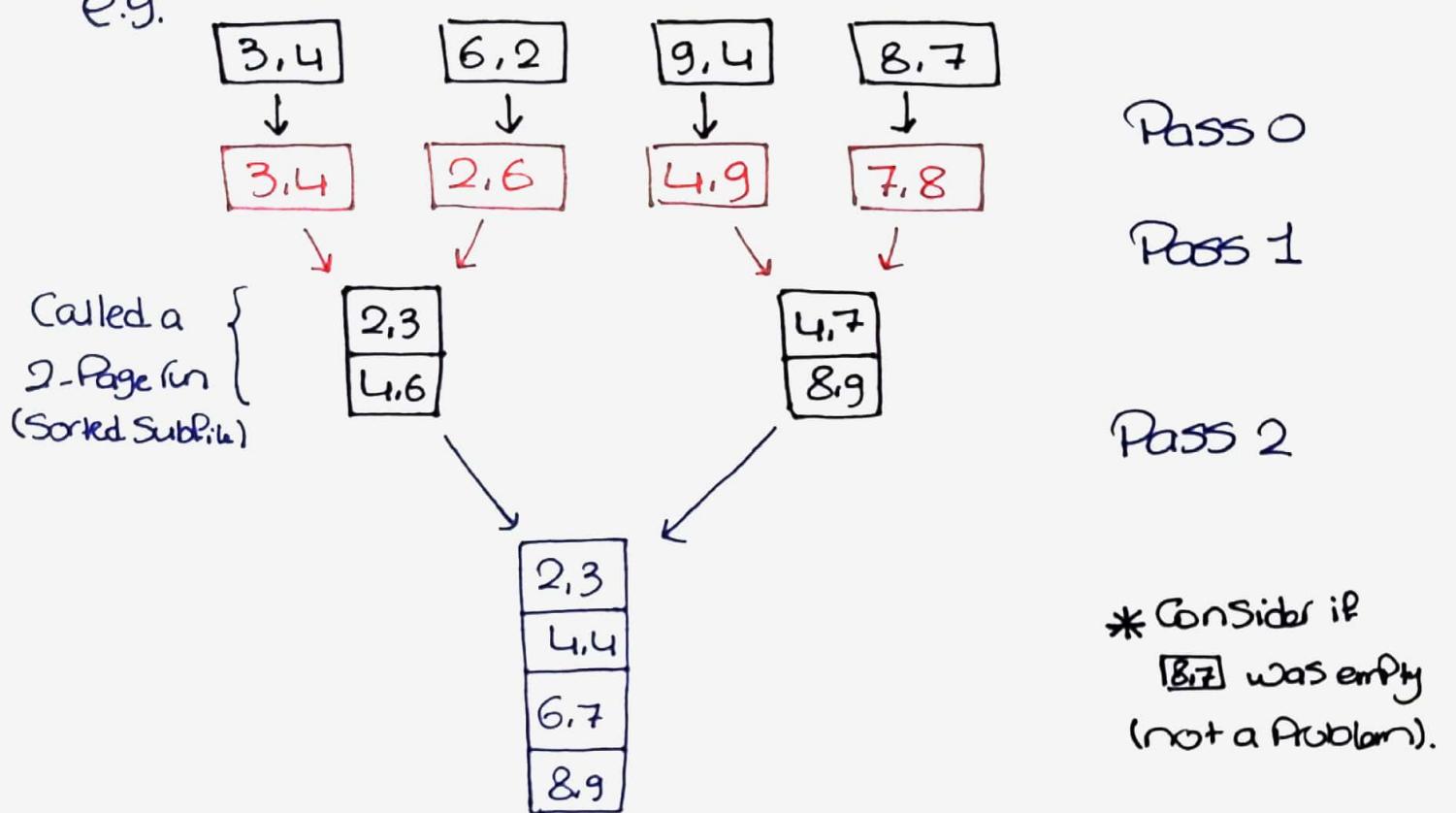
Only one
input buffer
used.

- Pass 0
 { read all 2^k pages
 & Sort each

- Pass 1
- read 2 Pages and Sort merge them into Output buffer
- Clearly need to do that 2^{k-1} times

- Keep going (Pass 2, 3, 4, ..., K) = Sorted File

e.g.



→ Zoom in buffers

[3,4] [2,6]

loaded into the two input

the two min. elements

[3,4]

[2,6]

→ [2,3]

write it

. then 4,6 remains,
write it next.

two reads and two writes

→ Zoom in

①
[2,3]
[4,6]

and

②
[4,7]
[8,9]

Decision:
load the next block
From the run ①
→ we know x > 7 for ②

[2,3]

→

[2,3]

write it!

[4,6]

→

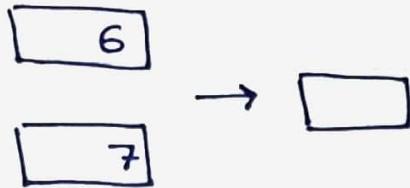
[4,4]

write it!

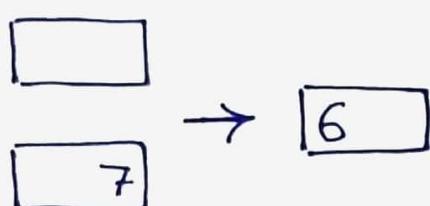
[4,7]

. read 1st two blocks
of each

Now



- Since $x > 6$ and $x > 7$
For ①, ② resp., Put 6 in OutPut



- load 8.9 (Final block)



—

→ Consider another SubProblem (to further clarify)

- Next block in each of the following runs



- buffers State was

From ① 1,5



→ 5

From ② 2,22



- now it is

* which run should we load from the next block?

Choose ① (now $x > 11$ for the next decision)

①

$x > 5$

②

$x > 22$

- The no. of Passes is clearly $\lceil \log_2 N \rceil + 1$ $(2^k \rightarrow N)$
 → in each, we read and write N blocks

• hence, # IOs = $2N(\lceil \log_2 N \rceil + 1)$

⇒ Suppose we have B buffer pages in main memory ($B-1$ as input and 1 as output)

Pass 0

- not the same as 2-way } → load the N blocks B by B and Sort internally InPlace
- result is hence $\lceil N/B \rceil$ sorted runs.

Pass 1, 2, ...

- like 2-way }
- Need the Output page (hence $B-1$) } → merge each $B-1$ runs (Sort-merge)

⇒ Initially have $\lceil N/B \rceil$ runs and tree branching factor is $B-1$, thus

$$\# \text{IO} = 2N(1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil)$$

↓ • Pass 0 • Pass 1, 2, ...

• read & write N blocks each

Pass.

- Example) Consider $B=5$ and $N=108$
 \hookrightarrow 4-way merge

Pass 0 → Pass 1 → Pass 2 → Pass 3

- Product $108/5 = 22$ • $22/4 = 6$ • $6/4 = 2$ • Coverage
Sorted runs Sorted runs Sorted runs (2-way merge)

• Works Surprisingly Well

→ $N = 10^9$ and $B = 257$ means we need
 $1 + \lceil \log_{256} 1 \rceil = 4$ Passes (excluding Pass 0)

→ although high B helps $\downarrow \#IO$ a lot, it also increases CPU time

• Next Slide (algorithm modifications) was Probably Skipped, has a lengthy discussion in the book and is Per Se Confusing.
(B-b input buffers, b output...)

- Consider when an input block becomes empty
 - need to suspend execution until its loaded again (then find min. over all block's...)
 - use an additional B shadow block's that prefetch from disk (B-1 and output)
 - Once input buffer becomes empty switch to a full shadow block and prefetch in the mean time.
- * CPU hence, is now never idle

6

Double Buffering

• Query Evaluation

- Needs to be optimized
- Two main issues

what are the possible query evaluation plans?

- * Ideally ← want best plan!
- * Practically ← avoid worst plans!

How is the cost of each estimated?

- use search algo. to find plan of best cost.

- In the following, we'll ignore the final write by the operation (as if its output will be used by a parent operation in the query tree; **Stays in memory instead of being written into disk**)

* Algorithms for Joins

Select *

From Reserve R, Sales S

Where R.Sid = S.Sid

} Need to Compute
 $R \bowtie_c S$

- Since $R \times S$ is large, it's inefficient to do that followed by Selection to implement the join.
- Want to optimize # IOs (excluding final write)

. Algorithms to consider

- Nested Loops Join
 - Simple and block versions

→ Sort Merge Join

→ Hash join

1. Simple Nested Loop Join

For each tuple r in R

For each tuple s where $r == s$

add $\langle r, s \rangle$ to result

- let M, P_R be the no. of blocks for R and the no. of tuples in each, respectively.
- N, P_S for S .

• then $\# \text{IOs} = M + (P_R \cdot M)N$

\downarrow $\underbrace{\quad}_{\text{For each tuple in } R}$ \downarrow
 R is scanned once
 (Outer loop) For each tuple in R Scan S
 (look for a match)

→ Simple refinement (Page-at-time)

(Page block)

for each Page m in R
 for each Page n in S
 add any $r \in m, s \in n, r=s$ to result } memory

$\# \text{IOs} = M + M \cdot N$

\downarrow \downarrow
 R is scanned once
 (Outer loop) S is scanned for each block in R

• Clearly, making the Outer loop the smaller relation results in less I/Os

• This uses 3 buffers (one to read R block another to read S block and third to store the output).

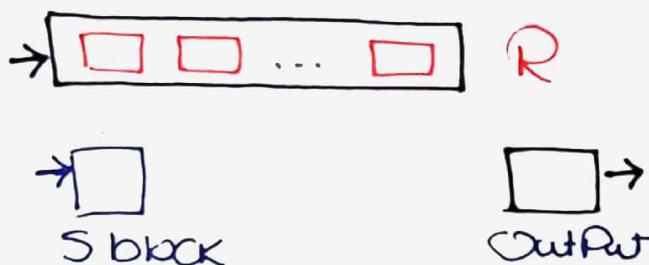
2. Block Nested loop Join

→ Suppose we have enough buffer pages

- to load entire smaller relation R
- Still have 2 extra buffer pages

- In this case, once we load in the smaller relation R (M IOs) we only need to read S block by block and write the output (2 extra buffers)

$$\# \text{ IO} = M + N$$



- Same for bop but R is all in memory

- What if the smaller relation R doesn't fit in memory?

→ Have B blocks
 → Load R Part by Part in $B-2$ blocks
 → 2 Blocks left for S and Output

Part of
 • each time $\rightarrow R$ is loaded, find all matching pairs for each block in S .

Clearly, R will be loaded over $\lceil \frac{M}{B-2} \rceil$ chunks

$$\rightarrow \# \text{ IO} = M + \lceil \frac{M}{B-2} \rceil \cdot N$$

The general nested loop join for blocks is hence

For each chunk of $B-2$ pages of R

For each page of S

• add $\langle r, s \rangle$ to result for all matches of the chunk and S block

3. Sort-merge Join

→ Sort each of R and S on the join Col.
(R.Sid and S.sid respectively) # external sort

racing { → Keep Scanning R until you find a tuple >,
the current tuple in S
→ Switch to S and do the same



- Whenever $r = s$, Output the qualifying tuples (tuples having the same value are next to one another)
- Start racing again until $r = s$ one more or both r and s are at the end of table.

S.Sid
→ 22
① 28 > 22 → 28
④ 31 > 30 → 30
58 > 44 → 44
58 > 44 → 58 . Match!
⑥

R.Sid
→ 28
28
→ 31
31
31
44 > 31 → 31
→ 58

update the 2 P+R ③
Formula we derived earlier (2NT)

. Cost:

$$\# \text{ IOs} = O(M \log M) + O(N \log N) + (M+N)$$

Sort R Sort S Traversal

- An extreme corner case would be that all of Sid in both relations have the same value, in this case, if the buffer pool is limited we may need to read all S for each block of R → MN instead of M+N
- If at least R or S has distinct values then M+N is guaranteed.

- Sort-Merge or block Nested ($M=500, N=1000$)

$B =$	35	100	300
# IOs (block n.)	16500	6000	2500
# IOs (sort n.)	7500	7500	7500

- More buffers make block nested better
- Results are by the previously shown formulae

3. Hash Joins

Partition Phase {

- Partition each of R, S into K Partitions on the disk using the same hash function h_i . (each Partition can span multiple blocks)

Probing Phase {

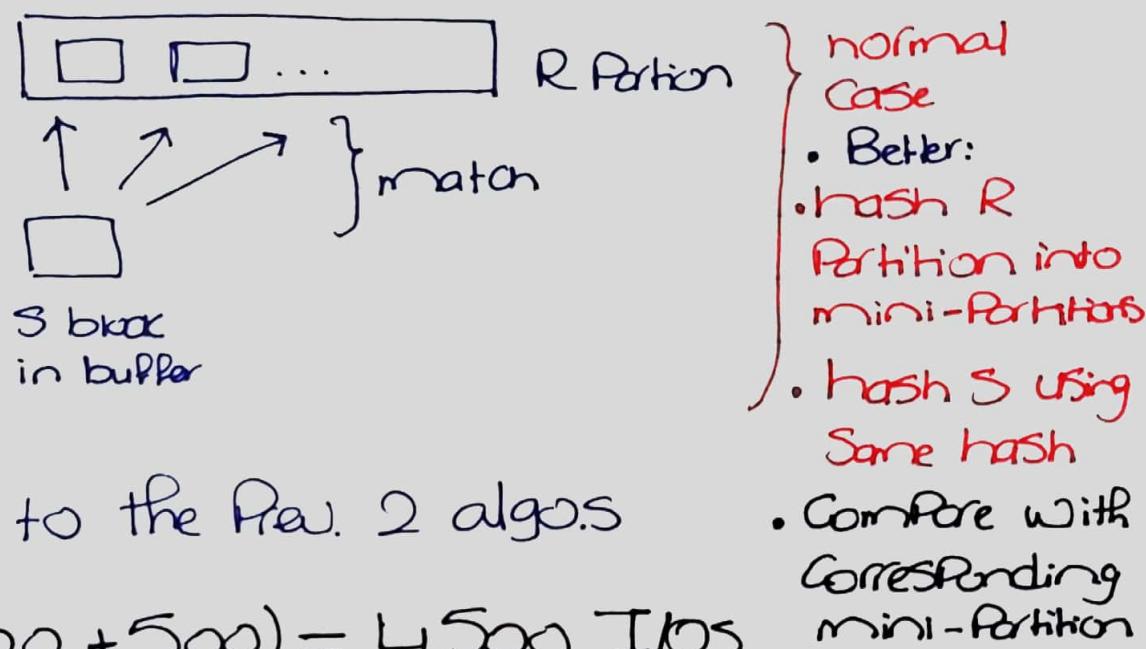
- Assume we can hold an entire Partition of R in memory (have enough buffers)
- then load R , Partition by Partition
- For each load the corresponding Partition of S (block by block is okay here)
- Check for matches!

* If $x=y$ then $h(x)=h(y)$ So once a Partition of R or a block of S is loaded, we never have to do that again in the future.

$$\# \text{ IOs} = 3(M+N)$$

- $2(M+N)$ in Partition Phase
- $M+N$ in Probing Phase

- To reduce matching effort by CPU
- We can use an in-memory hash table h_2
- hash entire R Partition and S block from Corresponding Partition by h_2
- Use it to match (much easier now)



Hash Join vs. Sort-Merge

$3(M+N)$
(also memory assumption)

- highly Parallelizable (hashing)

$3(M+N)$
if $B > \sqrt{N}$ (has Proof)

- less sensitive to data skew & gives sorted result.

* We implicitly assumed in hash join that each relation was ideally uniformly hashed (no data skew)

so its just as if we resorted the file (no extra space)

* Algorithms For Selection

Select *
From R
Where R.name=...

- No index, unsorted
 - Scan entire relation
 - expensive/inefficient if low selectivity
- No index, Sorted
 - binary search to find 1st tuple then read all matches
- B+ Tree or Hash index
 - consider cost of accessing data entries (leaves of B+) and matching
- Complex Conditions (day < ... and bid = ... and ...)
 - use one index then filter (e.g. on one of the keys)
 - use two indexes (on day and bid), intersect then apply 3rd condition (filter)

* Algorithms For Projection

→ removing fields is easy

→ removing duplicates (DISTINCT specified) isn't

- Sort-based
 - Sort to eliminate duplicates
 - remove unwanted attributes before sort
 - more standard than hash ← memory assumption
- hash-based
 - hash-join (relation with itself to find duplicates)
- index-only scan
 - attributes ⊆ index attr.

* Algorithms for Set Operations

- Intersection, Cross Products are a Special Case of Joins
- For union, Except: Sort and hash options
Very Similar to Joins and Projection.

* Algorithms for Aggregate Operations

- Sum, avg., min (Similar to Proj. approaches)



- Grouping
 - Sort or hash on group by attribute (or both!)
 - index only Plan can be Possible



- No Grouping
 - generally, need to Scan it all in this case
 - if index has all attrs in Select/WHERE then index only Plan is okay