

**ORDB**

Type T

DB

built-in  
integer  
char  
Number(10,2)  
varchar2

UDT

attributes  
methods

table

$T_1, T_2, T_3$

Referencing

Create type student\_t as object

sno  
sname  
gpa

Course ref Course\_t

Create type course\_t as object

Inserting data referenced

Insert into course values (c01)

Insert into student values (student\_t)

(select ref(c) from course c  
where c.cno = 'c01')

Select with references

Select with obj

Create table student

Create table student of student\_t

Select sname, scurrent  
from student s

Select & from student

ref returns the object id of the  
obj

Insert into student values (student\_t)  
(c01, 'Samon', 3.7);

Selecting as object  
as record

Select & from student

Select values (c)  
from student s

**VARRAYS & NESTED TABLES**

S001	kamal	D01	F01	S01	kamal	I-1.1
S002	Bimal	D02	F02	S02	Bimal	III
		D03	F03			

Create type price\_array as varray(10) of  
numbers(10,2)

Create table price\_list (

Pro int  
Prices price\_array )

Insert into price\_list values (c01, price\_array C  
250, 375, ...)

2 types of objs in tables

row  
Column

row	row
1	2

Column

Pro	Prices
1	2

② Selection from a Varay renaming.

Select Proj, 5. Column - Value Prices  
from pricelist P, table(p.prices) 5;

Noted tables

Employee

create type project for object  
proj --  
proj --

create type proj-list as table of proj  
create table emp as object  
type  
eno number  
projects proj-list)

table creation

Create table 'emp' of emp-t(eno  
nested table projects store as employee  
renaming)

insertion

Insert into employees values (1, Project-list (proj( ), proj( ) ..) ..)

methods

create type menu-t as object  
begin  
price float  
member function price in Var(Crate float)  
returns float  
)

create type body project menu-t as  
member function  
begin  
return rate \* self.price  
end  
end

add another member function

alter type menu-t  
add member function price in USD (rate float)  
return float;  
cascade.

create or replace body type menu-t  
if price in yes implemented  
member function price in USD (rate float)  
begin  
return self.rate  
end  
end.

Object / conversion

map / order

create type rect-t as object  
(length number,  
width number);

Map member function area returns number;  
)

create type body rect-t as object  
map member function area returns number  
begin  
return self.width \* self.height;  
end  
end

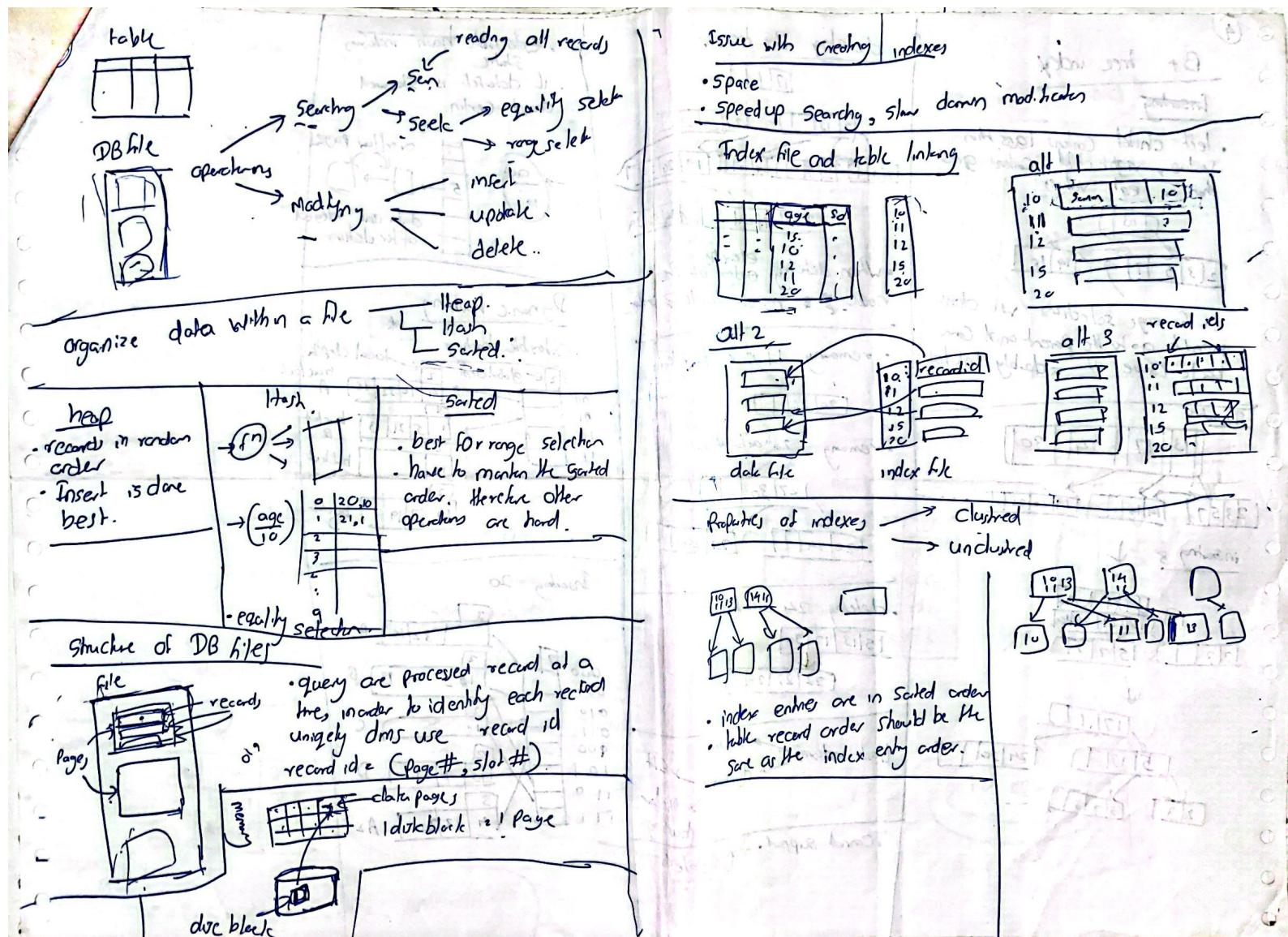
select distinct value(r) from rectangle r;

select value(r)  
from rectangle r  
order by value(r)

order by

create type curv-t as object  
id number  
name varchar(20)  
order member function match(c curv-type)  
returns integer

create type body curv-t as object  
order member function match(c curv-type)  
returns integer is  
begin  
if id < c.id then return;-1  
else if id > c.id then return;1  
else return 0;  
end

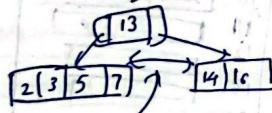


c (4)

### B+ tree index

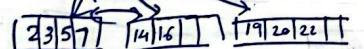
#### Inserting

- left child contains less than value, right child contains greater than or equal values

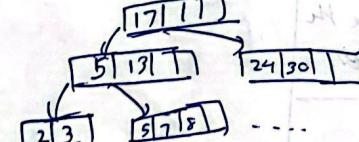
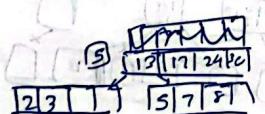


for range selection we don't need to go to the parent and come back cause of addressable link ind.

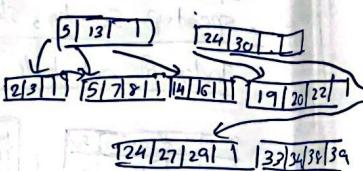
#### Deletes



#### inserting 8

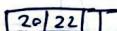


### deleting B+ tree

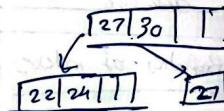


- When deleting element order of the node  $\geq 3$ , no of elements  $\geq$  order.

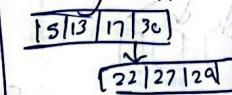
- removing 19 is no problem



- removing 20 after that



- delete 24 after that



- Hash indexes  $\rightarrow$  static  $\rightarrow$  dynamic  $\rightarrow$  EHT  $\rightarrow$  linear hashing
- Cannot support range selection

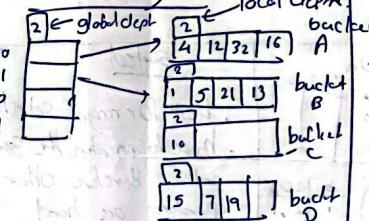
### Extensible hash indexes

- if dataset is skewed
- Grand Cycles

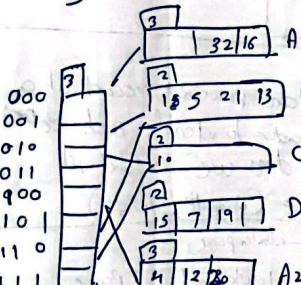


### Dynamic hashing

#### Extensible hash



#### Inserting 20



relational algebra

Input Relation - 1 or more  
Output relation - 1

relation : Set of tuples. set theory  
 $\{S, 8, 10, 11, 12, 18\}$

$R_1$		$R_2$	
$p$	$a$	$q$	$r$
$p_1$	$q_1$	$q_1$	$r_1$
$p_2$	$q_2$	$q_2$	$r_2$

Natural join (\*)

$R_1 \bowtie R_2$

$R$	$p$	$q$	$r$	$s$
$p_1$	$q_1$	$r_1$		
$p_2$	$q_2$	$r_2$		

equi join ( $\bowtie$ )

$R_1 \bowtie R_2$   
 $R_1.p = R_2.s$

$p$	$q$	$r$	$s$
$p_1$	$q_1$	$r_1$	$s_1$
$p_2$	$q_2$	$r_2$	$s_2$

$S$	$R$
$q_2$	$r_2$
$q_3$	$r_3$

$P$	$a$
$p_1$	$q_1$
$p_2$	$q_2$

Project operator ( $\pi$ )

$\pi(sno)$  (Students)  
 $sno, gpa$

sno	name
1	Kunal
2	Nimol
3	Ishan

sno	sname	gpa
1	Kunal	3.8
2	Nimol	2.5
3	Ishan	3.7

join

Relational Algebra

Renaming relation (P)

$P_{T_1} \leftarrow \pi(sno)$   
renaming tables

$P_{student}(sname, sno) \leftarrow \pi(sname, sno)$   
renaming fields.

Select sno as student no  
from student - S  
where gpa > 3.0  
or  
select sno [student no]

Select operator ( $\delta$ )

( $\delta$  (Students))  
 $gpa > 3.0$

select \*  
from student  
where gpa > 3.0

sno	sname	gpa
1	Kunal	3.8
3	Ishan	3.7

operators

AND, OR, NOT

Cartesian product ( $\times$ )

A	B	C	D
$a_1$	$b_1$	$c_1$	$d_1$
$a_1$	$b_1$	$c_2$	$d_2$
$a_2$	$b_2$	$c_1$	$d_1$
$a_2$	$b_2$	$c_2$	$d_2$

Select  
From student

Select \*  
From R, S.

Heuristic join ( $\bowtie$ )

$R_1 \bowtie R_2$   
 $R_1.p \neq R_2.s$

not equals

union (U)

R		S	
A	B	C	D
a <sub>1</sub>	b <sub>1</sub>	a <sub>2</sub>	b <sub>2</sub>
a <sub>2</sub>	b <sub>2</sub>	a <sub>3</sub>	b <sub>3</sub>
a <sub>3</sub>	b <sub>3</sub>	a <sub>4</sub>	b <sub>4</sub>

RUSN →

A	B	C	D
a <sub>1</sub>	b <sub>1</sub>	c <sub>1</sub>	d <sub>1</sub>
a <sub>2</sub>	b <sub>2</sub>	c <sub>2</sub>	d <sub>2</sub>
a <sub>3</sub>	b <sub>3</sub>	c <sub>3</sub>	d <sub>3</sub>
a <sub>4</sub>	b <sub>4</sub>	c <sub>4</sub>	d <sub>4</sub>

RUS →

A	B
a <sub>1</sub>	b <sub>1</sub>
a <sub>2</sub>	b <sub>2</sub>
a <sub>3</sub>	b <sub>3</sub>
a <sub>4</sub>	b <sub>4</sub>

Intersection (n)

R ∩ S →

A	B
a <sub>2</sub>	b <sub>2</sub>
a <sub>3</sub>	b <sub>3</sub>

Left outer join (IX)

Emp Dept

eno	ename	deptno	dname	loc
e <sub>1</sub>	Sam	D <sub>1</sub>	SE	PT
e <sub>2</sub>	Konal	D <sub>2</sub>	IT	NULL
e <sub>3</sub>	Nimal	NULL	NULL	NULL

Full outer join (IX)

eno	ename	deptno	dname	loc
e <sub>1</sub>	Sam	D <sub>1</sub>	SE	PT
e <sub>2</sub>	Konal	D <sub>2</sub>	IT	NULL
e <sub>3</sub>	Nimal	NULL	NULL	NULL
NULL	NULL	D <sub>3</sub>	DS	NULL

Relational Algebra

grouping (J)

aggregate funcs

eno	ename	sal	dept
1	A	15000	D <sub>1</sub>
2	B	15000	D <sub>1</sub>
3	C	30000	D <sub>2</sub>
4	D	10000	D <sub>2</sub>

P(T<sub>2</sub>(dept, count, avg(sal)))

ex3:

J(dept, count(\*), avg(sal))

dept	Count	Sal
D <sub>1</sub>	2	12500
D <sub>2</sub>	2	20000

ex1: J(count(\*))

[4]

department wise emp count

ex2: P(T<sub>2</sub>(dept, emp count))

J(count(\*))

dept	Count
D <sub>1</sub>	2
D <sub>2</sub>	2

## query Processing

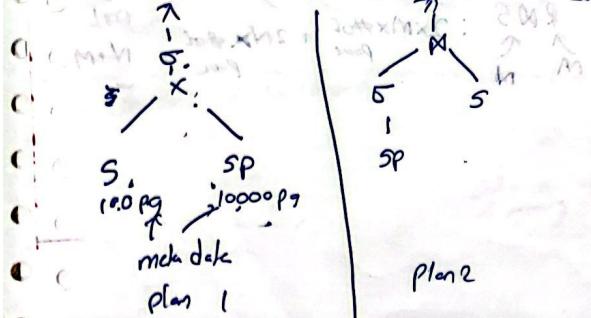
- ① Parsing and translation
- ② Optimization
- ③ Evaluation.

Select S.name  
from S, SP  
Where S.Sno = SP.Sno AND  
SP.Pro = 'P<sub>2</sub>'.

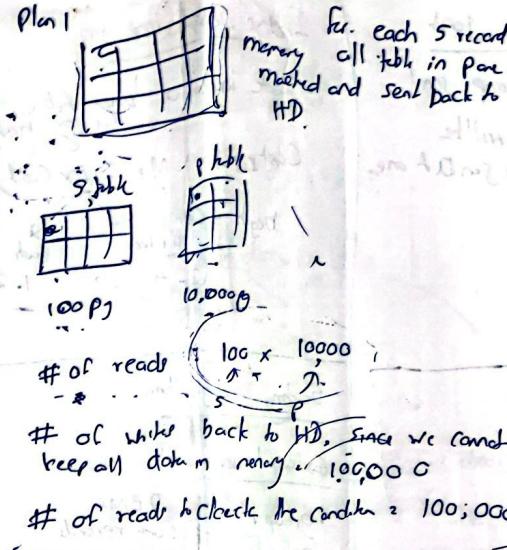
$\pi_{Sname} (\sigma_{SP.Pro = 'P_2'} (S \bowtie SP))$   
Since S.Sno = SP.Sno  
AND SP.Pro = 'P<sub>2</sub>'.

$\pi_{Sname} (S \bowtie \sigma_{SP.Pro = 'P_2'} SP)$

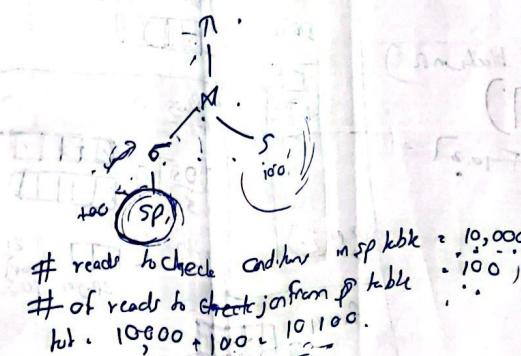
To find the best execution plan, we calculate cost: 10 operations.



Plan 1



Plan 2

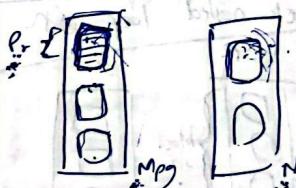


## Algorithms for joins

- Simple nested loop join
- Page oriented Nested loop join
- Block nested loop join
- Index nested loop join
- Sort merge join

### Simple nested loop join

R  $\bowtie$  S



Cost:  $SNLJ = M + (MPr \times N)$   
since we take one page of R for each rec, inner table is scanned.  
for each tuple in outer table the inner table is scanned.

⑧

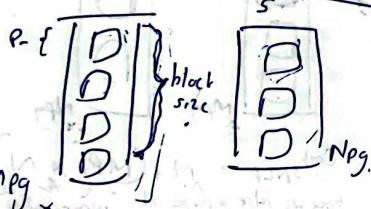
### Page oriented nested loop join

- Unlike simple nested loop join!
- In here, 1 page of supplier will be compared to the entire table to select one row.

Selected

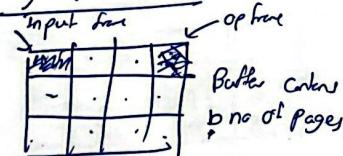
$$\text{Cost}_{PONLT} = M + MN.$$

### Block nested loop join



$$\begin{aligned}\text{Cost}_{BNLJ} &= M + N \times (\# \text{ of blocks in } R) \\ &= M + N \left( \lceil \frac{M}{b} \rceil \right)\end{aligned}$$

Finding # of blocks



### Index nested loop join

- Create index of inner table  $B_1$  or hash.

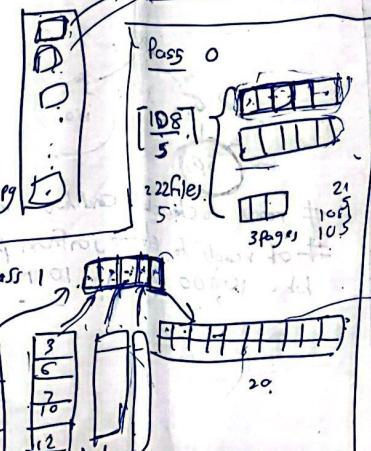
$$\text{Cost}_{INLJ} = M + (\text{Index Cost}) (M \times P)$$

$$\begin{aligned}\text{height} &= \log_{B_1} \text{tree} \\ &= (2-4) \text{ bits} \\ &= 1.2\end{aligned}$$

### Sorting algorithms

#### External merge sort

- group by DBMS
- order by Sort records
- disk
- store in buffer
- dump



### Pass 0

$$\left\lceil \frac{108}{5} \right\rceil = 22 \text{ files}$$

5 pages each file  
 $21 \times 5 + 3$

### Pass 1

$$\left\lceil \frac{22}{4} \right\rceil = 6 \text{ files}$$

20 pages each file  
 $20 \times 5 + 8$

### Pass 2

$$\left\lceil \frac{6}{4} \right\rceil = 2 \text{ files}$$

8 pages each file  
 $\frac{2}{3} = 1$

$$80 \times 1 + 28$$

- Cost: no of reads and writes
- $2 \times \# \text{ of pages} \times \# \text{ of pages}$
- $2 \times 4 \times 108$

### Sort merge join algorithm

$$\text{Cost}_{SMJ} = \text{Sizing Cost} + \text{Merge Cost}$$

$$= \text{Sizing Cost}(R) + \text{Sizing Cost}(S) + \text{Merge Cost}$$

$$\text{RNS} = 2 \times M \times \# \text{ of pages} + 2N \times \# \text{ of pages} + N + M$$

## Physical DB design & DB tuning

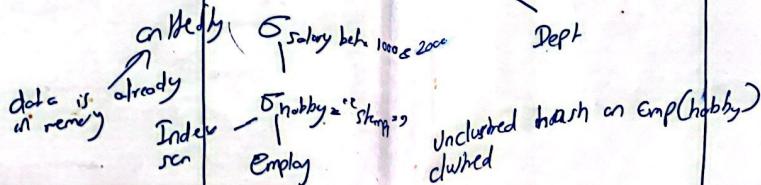
- Index selection & creation
- Tuning the conceptual schema
  - alternative normalized schema
  - Denormalization
  - Vertical partitioning
- Query rewrite.

## When selecting an index

- What relations do I need
- What are the search fields
- Single attributes or multi-attribute
- type? clustered / unclustered
  - B+ tree or hashed
- Impact on update
  - hash indexes — equality selection
  - B+ tree — range selection.

Most selective one is taken first in a multi-attribute index.

clustered v unclustered.



## Heuristic rules

- Avoid cartesian product
- Bring selection down
- left d plan

## Tuning Conceptual schema

- denormalization
- Vertical partitioning
- Views

- DB programming level
- Application Coding level

- stored procedure
- begin T
- commit T
- connection object
- Set properties
- Combo Commit file
- Set Combing pub

## Isolation level

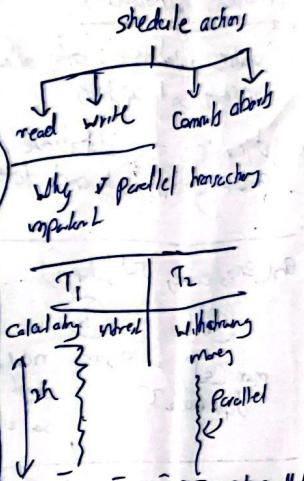
- Serializable (S2PL + index locks)

- Repeatable Read (R2PL)  
Phantoms possible

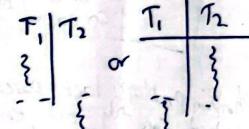
### Transactions and Concurrency Control

- Transaction: Sequence of reads & writes that belongs to one unit of work.
- Benefits from the database side / properties
  - A - Atomicity
  - C - Consistency
  - I - Isolation
  - D - Durability
- Atomicity: all actions of the transaction happens or nothing happens. (Executed as a single unit of work)
- Consistency: Maintains integrity of dataset in the DB (Group of constraints)
  - key
  - table
  - referential
  - domain
  - assertions
- Isolation: Even though actions of several transactions are interleaved, the net effect is identical to executing all transactions in some serial order.
- Durability: After transaction is committed, the data should be persistent.

(Transactions survives crashes) & system failures



Serial schedule: Schedule does not interfere the actions of multiple transactions.



- no conflicts
- always ends with correct schedules.

equivalent schedule: same effects from the schedules despite execution order. For any db state, effect of executing serializable schedule: (and T1 & T2) things are happening parallelly but equivalent to some serial schedule.

### Anomalies in interleaved transactions

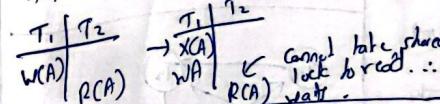
- a) W(A), R(A) (dirty reads 1.02.21.)
- b) R(A), W(A) (non-repeatable reads)
- c) W(A), W(A) (overwriting unwanted) data

### Lock based Concurrency protocol

#### Strict 2 phase locking (Strict 2PL)

- transactions must obtain a shared lock before reading and exclusive lock before modifying
- All locks for a given transaction is released after committing (Complete)
- If a transaction has an exclusive lock on an object, other transactions cannot have a shared or exclusive lock on it.

#### How 2PL solves cascading

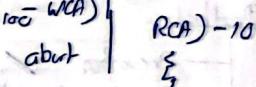


- Cascading aborts
- Unrecoverable abts.

### Cascading

$$A = 50$$

$$T_1 \quad | \quad T_2$$



### How 2PL solves cascading

$$T_1 \quad | \quad T_2$$

$$\uparrow X(A) \quad \uparrow W(A)$$

$$\downarrow R(A) \quad \downarrow C(A)$$

$$\downarrow A \quad \downarrow C$$

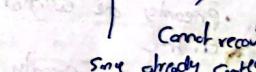
### Unrecoverable

$$T_1 \quad | \quad T_2$$

$$\uparrow W(A) \quad \uparrow R(A)$$

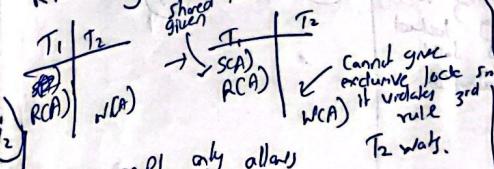
$$\downarrow R(A) \quad \downarrow C(A)$$

$$\downarrow A \quad \downarrow C$$



Cannot recover since already Committed

### RW Slicing



store/recover 2PL only allows serializable schedules.

### overhead of lock based CC

- Monitoring & handling of locks
- Handling deadlocks
- Handling phantom.

### Deadlock

- Prevention ✓
- Avoidance (detection)

### Cycle of wait.

### Prevention methods

- + wait & die
- + wands & wait

### Wait & die



SC(A)  
RC(A)

X(C(B)) wait  
for S(C(B))

X(C(A)) wait  
for S(C(A))

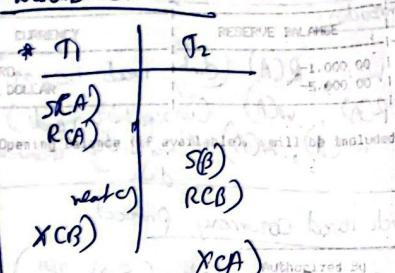
- Priorities are assigned based on tree

T<sub>1</sub> is high priority .. it waits

But when T<sub>2</sub> is looking for X(C(A)), it is already locked

- and since it is low priority , it is aborted .

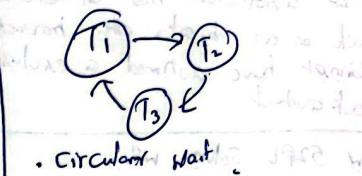
### Wand & wait



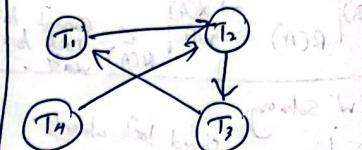
- High priority doesn't wait for low priority
- Wand the low priority and cancels it and take resource

### Detection

- Create a wait for graph



- Circular Wait



### phantom problem

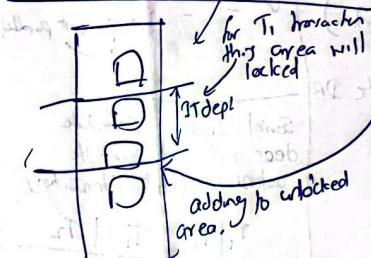
DATE : 10-06-2021  
5:13:00

static database: only allows inserts and updates

dynamic : database is growing insert and delete

Ex(eno, enat, sal, dep)

(T<sub>1</sub>)  
max (sal) = 100,000,  
dept <= 55  
new emp  
dept IT  
sal < 160,000



Phantom Problem: lock based concurrency protocol can only lock existing data in the database but parallelly running other transactions can manipulate data into other areas conflicting with

### Solution: Predicate lock in

### Index lock in

- Create an index on the selection criteria C dependent in the above query entries along with lock the index

### Index lock in

- Problem? It locks all from root to leaf.
- since in SP2I locks are released after completing transactions

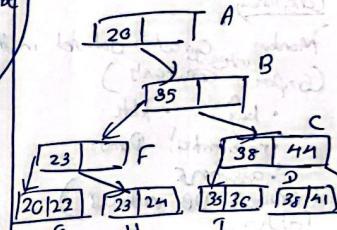
### Simple tree locking algorithm.

Search: Start from root & go down. Go to child and release parent.

Set

Insert: Node is not full

Delete: Node is not half empty



Delete 38

X(CA)

X(CB)

X(CC)

X(CD)

X(CE)

rel A, B

rel B, C

rel C

rel D

rel E

rel A

rel B

rel C

rel D

rel E

rel A

rel B

rel C

rel D

rel E

Lock management

Multiple granularity

Lock scheme

g. locks more efficient than (X)

table

How does this solve?  
Before releasing a lock, the lock manager wants to search, in many objects.

Solutions: multiple granularity levels

db  
↓  
table  
↓  
Page  
↓  
record.

Crash recovery

The way of guaranteeing atomicity and durability.

Why it is needed?

- before going to the commit part might write to disk
- committed transactions might not be written to disk

Concurrent data already gone to disk

- undo
- Force

at committing part, things are written to disk.

Redo

No Steal	Steal
Force	Trivial
No force	Desired

Steal

Can we write pages to disk of uncommitted transactions?

Data goes to disk from uncommitted transactions.

Write ahead logging protocol

- Must force log records for an update before corresponding data pages goes to disk (guarantees atomicity)
- Must write all records for transaction before commit (guarantees durability)

Log Sequence Number

LSN

log header

log file

Page

most recent modification

fluked LSN

Before a page is written: Page LSN

Log records

has:

- XID
- Type

Possible log record types:

- update
- insert
- commit
- end

transaction table

TID	Status	Last LSN
1	Active	100
2	Active	200
3	Active	300

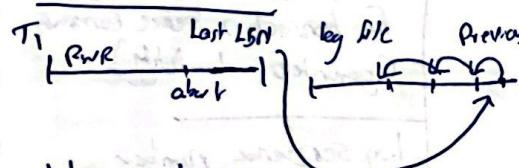
Dirty Page Table

Page #	Record LSN	data
1	100	old data
2	200	new data

### Checkpoints

- Creates checkpoints to reduce the time taken to recover in the event of a system crash.

### Simple transaction Abort



- Data page has the last LSN
- Along the log record which has the previous LSN we can go back.
- Before starting Undo, write an Abort log for recovering from crash during aborting.

To perform undo, it must have a lock on data.

- CLR entry is again written back to the log just in case a failure happens during the aborting process.

Last LSN (CLR)  
undone next LSN

Currently writing

prev LSN of the log = 245

- at the end of undo "end" is written.

### Commit

- Write Commit to log file after logging the record upto Xact's Last LSN are flushed.
- Guarantees that  $\text{flush LSN} \geq \text{last LSN}$

### ARTIES ACRES algorithm

- There are 3 phases in the algorithm

- Analysis to figure out what Xacts have failed since the last committed checkpoint.

- Redo all actions
- Undo failed Xacts

Smallest LSN      X-table      Last LSN  
in dirty page      oldest rec of      xact active at  
Tart checkpoint      table      crash



### Analysis Phase

- reconstruct the state at checkpoint via end-checkpoint record.
- Scan log forward from checkpoint.
  - End: remove these Xacts from the Xact table.

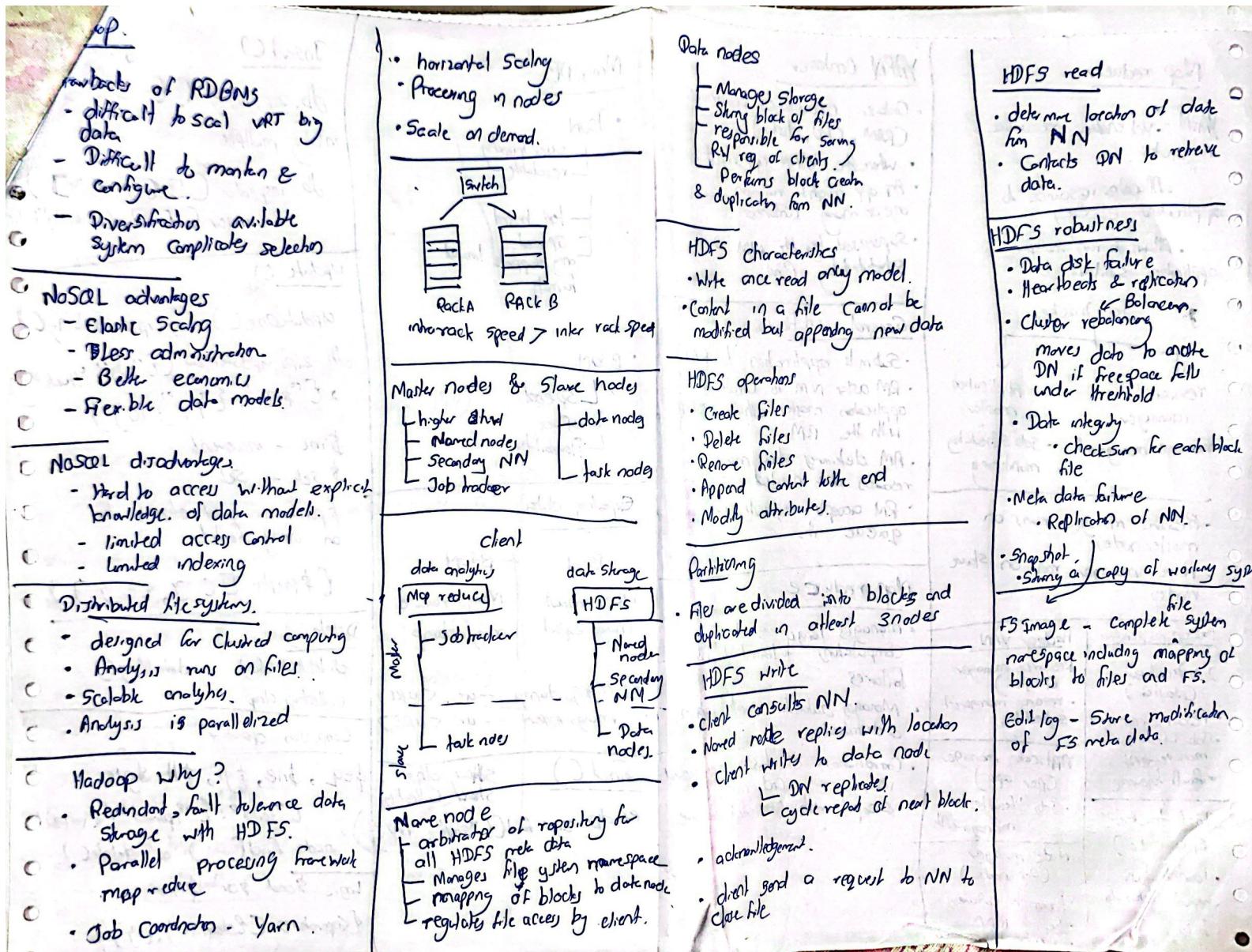
Analysis starts at the most recent

begin-checkpoint and proceeds forward in the log until the last log record. It determines:

- point in the log at which to start the redo phase
- the dirty pages in the buffer at the time of crash
- transactions that were active at the time of crash which needs to be undone

Redo phase: Redoes all changes to any page that might have been dirty at the time of crash.

Undo: Undoes all changes of all xacts that were active at the time of crash.



<u>Nop reducer</u>	
YARN - yet another resource negotiator.	
Allocates resources to applications effectively.	
Allows to run multiple applications simultaneously.	
<b>Job tracker</b>	
resource manager	Application master
• resource management	• job scheduling • monitoring
Resource manager runs on master nodes.	
Node manager runs on slave nodes.	
<b>Hadoop mapreduce</b>	<b>Hadoop YARN</b>
Job tracker (Master)	Resource manager (Master)
• resource management	• resource management
• Job lifecycle management.	• scheduling
• fault tolerance	
<b>Task tracker</b>	<b>Application manager</b>
(Per node)	(Per app)
• launch task	• job lifecycle management
• report status	• launch containers

<u>YARN Container</u>	
Container : Collection of resources (RAM, CPU, disks)	
where the YARN app runs.	
An app might run on one or many containers.	
Supervised by the NM & scheduled by RM.	
<b>General workflow</b>	
• Submits application to RM	
• RM asks NM to create an application master which registers with the RM.	
• AM determines how many resources are needed.	
• RM accepts requests and queues it.	
<b>Nop reduce</b>	
manages large scale computation tolerant to hw failures.	
• Manages automatic parallelization & distribution.	
• Coordination of tasks implemented in map() and reduce() & copes with unexpected failure.	

<u>MongoDB</u>	
JSON	
User friendly	
Readable	
Text based	
Spaced	
only supports limited formats	
<b>Insert()</b>	
db.zip.insert({})	
inserting multiple	
db.insert([{}], {check: false})	
<b>Update()</b>	
updateOne(), updateMany()	
db.zip.updateMany({city: "Helsinki"}, {\$inc: {"pop": 103}})	
\$inc - increment	
\$set - set	
\$push - add elements to an array field.	
{\$push: [{x: 3, y: 3}]}	
<b>Delete()</b>	
deleteOne(), deleteMany()	
Collection.drop()	
<b>Comparison operators</b>	
\$eq, \$ne, \$gt, \$lt, \$gte	
{"field": {"operator": < value}}	
inside find() or delete()	
<b>Logic operators</b>	
\$and \$or \$nor	
{operators: [{"statement": "x > 3"}, {"statement": "y < 5"}]}	

## Basic \$expr

`{$expr : {expression}}`

Compare fields diff user defined

`[$expr": {$eq": { "fieldID": "prodID", "prodID" }}`

## Array

`find ($c "category": {`

`"size": 20,`

`"fall": { "category": "winter",  
"children": "Hiking" } } } )`

## Projection

`find ($c .? { "price": 1, "address": 1, "id": 0 },`

## aggregation framework

`.aggregate ($math: "sum", "min:",  
{$proj: { "price": 1, "address": 1 } })`

## Sort() limit()

`find(). sort({ "pop": -1, "city": 1 })`

`db. hpi. createIndex ({ "buidYear": 1 })`

## Compound index()

`{ ... , ... }`