

Database Management Systems

Database Management System is basically a collection of data related to an enterprise with some programs to manage or manipulate it.

Program | Data

Isolation

Abstraction

(So that the low level details need not be known to the programmer)

- Convenience of Data
- Efficiency

Database Model

Relational Model

Every single unit of data is represented as a vector, where each component is called an attribute.

attribute

↓								
---	--	--	--	--	--	--	--	--

Data object / tuple / Record

Multiple records give rise to a "data table"

Eg.!

Table / Students of DBMS

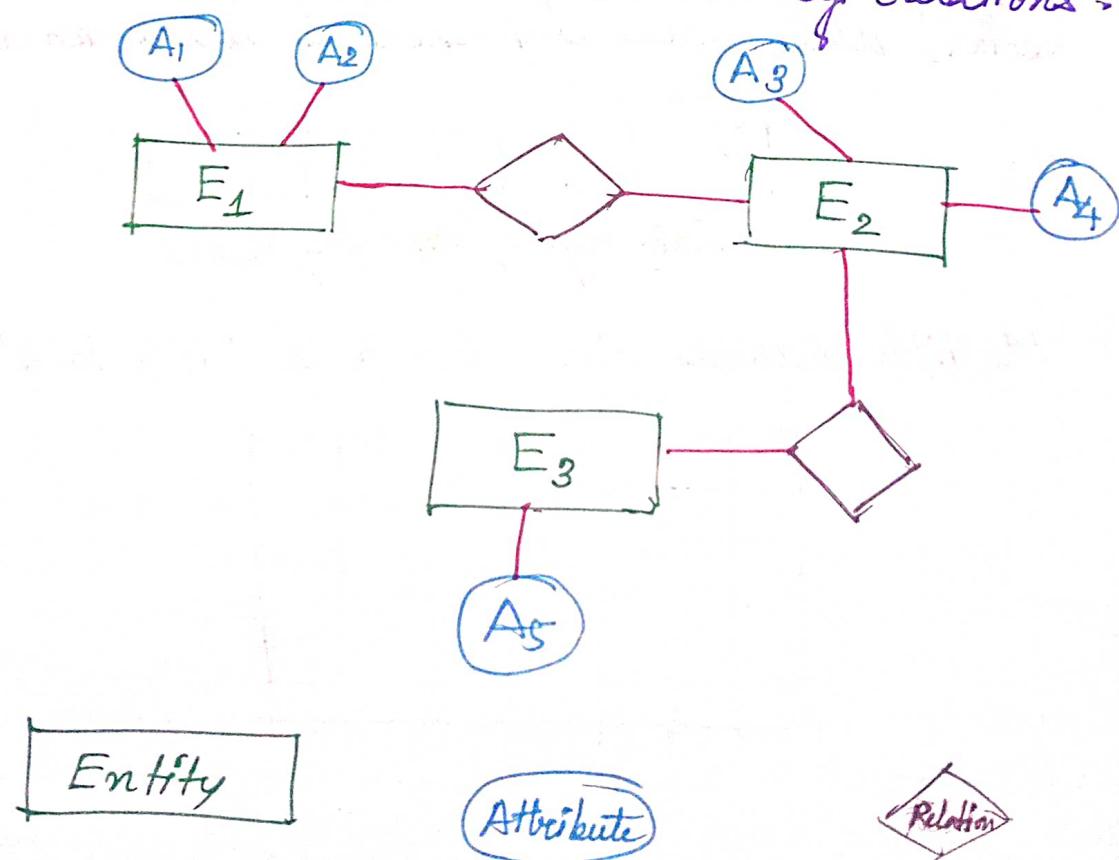
	RN	N	D	CG	SG	H
S ₁						
S ₂						
S ₃						
S ₄						

We can name this relation
"DBMS course students"

This is a subset of the cartesian product of everything.

- We will interchangeably use "Table" & "Relation"
- Structure of Record \rightarrow Schema of data ; All schema — Schema of DB
- Entity Relationship Model (ER- Model)

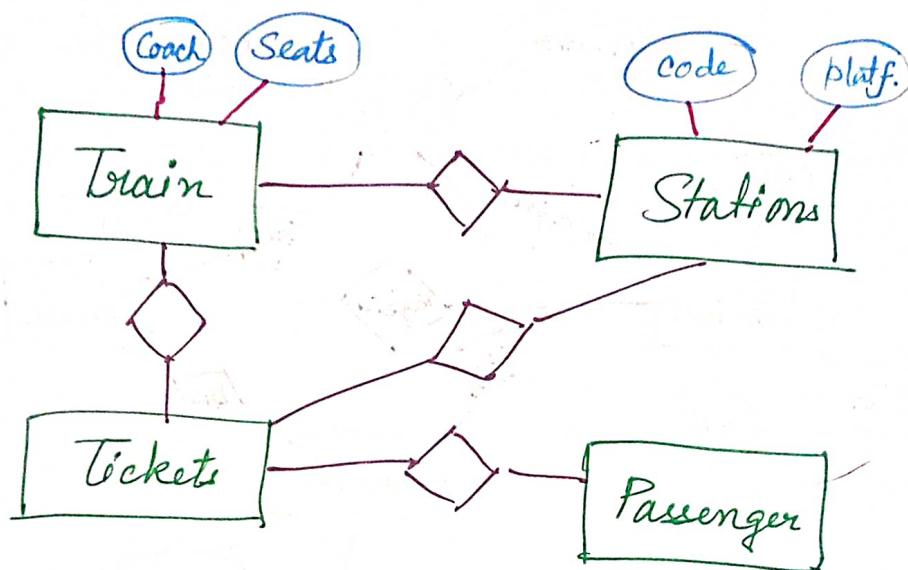
Data is represented in a diagrammatic way. There are entities (described by a set of attributes) and these entities are related by relations.



Eg.: ER-Diagram of IRCTC Booking System

Entities

- Train
- Passenger
- Stops/Stations
- Ticket



We can convert an ER-Diagram to a Relational Schema & vice-versa.

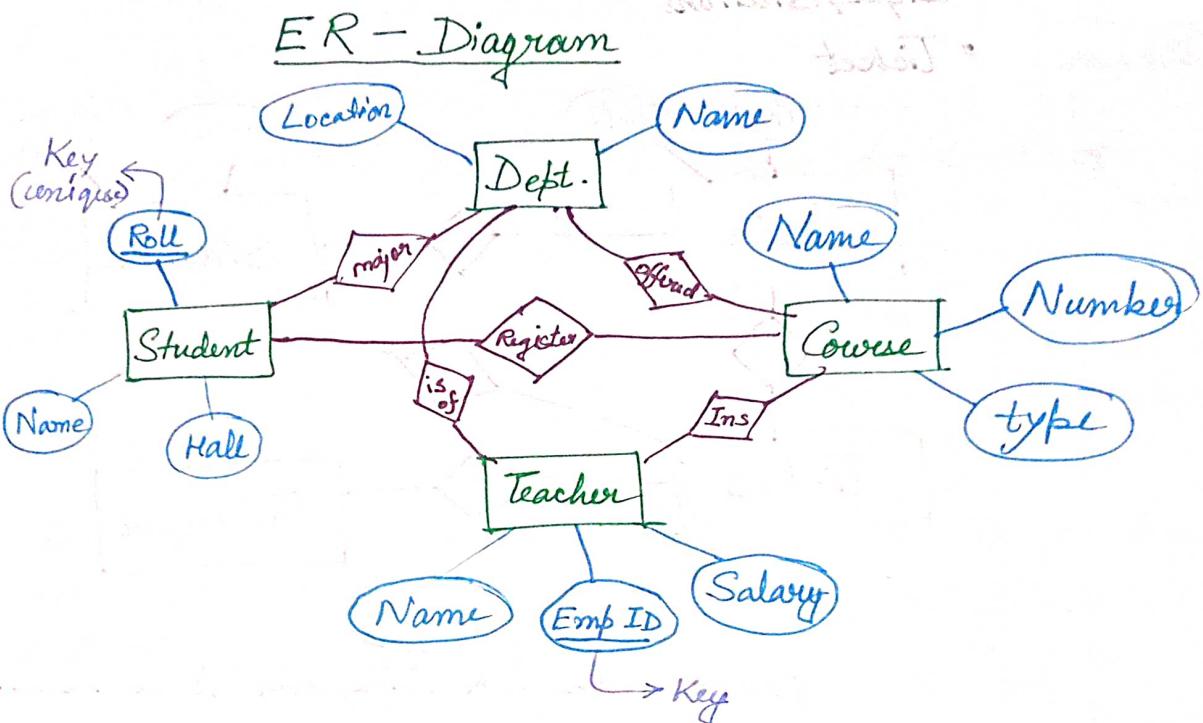
(We can have entities as tables & we have a common table entry to identify relation).

- ER table & Relational schema can be interrelated
- Relational schema cannot handle data without well defined set of attributes (unstructured data)
- Relational schema cannot handle data with varying number of attributes.

Semi Structured Data Model (XML)

Some markups or tags were used to describe data.

Used in E-commerce sites as Flipkart or Amazon



Step 1: Write down all entities and their attributes.

Step 2: Identify relations and incorporate them in the schema.

Key: Unique attribute of an entity (no two instances of this entity has the same key)

Superkey: Subset of attributes which has unique values for every instance of the entity

Minimal Superkey: Superkey whose no smaller subset is superkey (Candidate key)

- Superkey & Candidate Key are part of the schema
- There can be multiple Candidate Keys.
- Keys, when forced to be unique, are called "Primary Key" and is considered to be a part of the schema

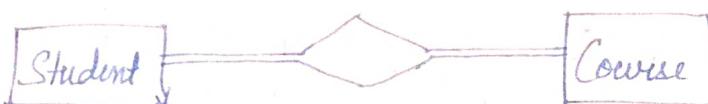
Primary Key: Candidate Key, which is forced to be unique.

Weak Entity Set: Entities which have no candidate key.

Weak Entity Sets
(double rectangle)

Participating Constraints:

(i) Total Participation: All entity instances of the entity have atleast one relation. (double line)



(ii) Cardinality

◦ Many to Many --

◦ One to Many ← -

◦ Many to One → -

◦ One to One: ← →



Converting ER-Diagram to a Table

	A_1	A_2	A_3	A_4	\dots	A_D
t_1						
t_2						
\vdots						
t_N						

table or

The table/relation schema is specified by

$$\sigma_L(A_1, A_2, \dots, A_D)$$

A database schema may be specified by -

$$\sigma_{L_1}(A_1, A_2, \dots, A_D)$$

$$\sigma_{L_2}(A_1, A_2, \dots, A_{D_2})$$

\vdots

Foreign Key:

σ_1	σ_2																																																		
<table border="1"><tr><td>A_2</td><td>A_5</td><td>A_7</td><td>A_8</td><td>A_9</td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr></table>	A_2	A_5	A_7	A_8	A_9																					<table border="1"><tr><td>A_1</td><td><u>A_2</u></td><td>A_3</td><td>\dots</td><td>A_D</td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr></table>	A_1	<u>A_2</u>	A_3	\dots	A_D																				
A_2	A_5	A_7	A_8	A_9																																															
A_1	<u>A_2</u>	A_3	\dots	A_D																																															

If A_2 is a key for σ_1 & σ_2 is a weak entity set, then A_2 is a foreign key referencing σ_2 from σ_1 .

E.g.:

Courses	
<u>C_id</u>	Cname

Student	
Roll	Cid

Cid is a foreign key referencing Student table from Courses table

Conversion from ER-Diagram to Table

- Have a table for
 - Each entity \rightarrow { Normal : All attributes
Weak : All attributes, plus attributes of all referencing relations of all referencing entities }
 - Each relation \rightarrow (key of both participating entities) (primary key)

Example:



RELATIONAL ALGEBRA

"Queries" are fired on table (or a set of tables). The answer may also be represented like a table.

result-table = Query(table1, table2, ...);

Can be specified by some Data Query Language.

Data Query Languages

- Procedural {Relational Algebra}
- Descriptive {Relational Calculus}

We start with Relational Algebra.

RELATIONAL ALGEBRA

Table \equiv Relation

α		
	A_1	A_2
t_1		
t_2		
t_3		

, $\alpha(A_1, A_2, A_3); t_1, t_2, t_3 \in \alpha$

This is a relation as, $\sigma_r(A_1, A_2, A_3) \subseteq A_1 \times A_2 \times A_3$

Now, we consider a relational algebra expression.

$$\sigma_r = (\sigma_{l_1} \phi_1 \sigma_{l_2} \phi_2 \sigma_{l_3})$$

In relational algebra, there are six primitive operators.

UNITARY OPERATORS

• "SELECT" operator

$$\text{result} = \sigma_\theta(\sigma_r), \quad \begin{cases} \text{Where } \theta \text{ is a clause, like} \\ \cdot (A_1 > 30) \& (A_2 = 0) \\ \cdot (A_3 == 0) \end{cases}$$

then, we have result as the set of tuples from σ_r satisfying θ .

$$\text{result} = \{t \in \sigma_r \mid t \text{ satisfies } \theta\}$$

• "PROJECT" operator

$$\text{result} = \pi_{\{A_{i_1}, \dots, A_{i_k}\}}(\sigma_r)$$

result will just be the table with only the columns of the attributes A_{i_1}, \dots, A_{i_k} . Mathematically, removal of duplicates are necessary.

In SQL, implementation does lead to duplicates.

- "RENAME" operator

$$\rho_{\sigma_1}(\sigma_2) ; \text{ renames } \sigma_2 \text{ as } \sigma_1$$

Example: student (Name, Roll, Hall, course).

Find the Name & Roll of students from RK, taking the DBMS course.

$$\sigma_{\text{res}} = \prod_{\{\text{Name, Roll}\}} \left(\overline{\sigma_{\text{student}}} \mid \begin{array}{l} \text{Hall} = 'RK' \\ \& \text{course} = 'DBMS' \end{array} \right)$$

BINARY OPERATORS

- "SET UNION" operator

$$\sigma_3 = \sigma_1 \cup \sigma_2 \quad \left\{ \begin{array}{l} \text{The attributes in } \sigma_1 \text{ and } \sigma_2 \\ \text{must be the same} \end{array} \right\}$$

Mathematically, σ_3 is without duplicates.

- To save computational power, SQL implements union by allowing duplicates.

- "SET DIFFERENCE" operator

$$\sigma_3 = \sigma_1 - \sigma_2 \quad \left\{ \begin{array}{l} \text{Again, the attributes in } \sigma_1 \text{ & } \sigma_2 \\ \text{must be the same} \end{array} \right\}$$

- "CARTESIAN PRODUCT" operator

$$\sigma_3 = \sigma_1 \times \sigma_2$$

if $\sigma_1(A_1, A_2, A_3)$ and $\sigma_2(A_4, A_5)$; $\sigma_3 = (A_1, A_2, A_3, A_4, A_5)$

if $\sigma_1(A_1, A_2, A_3)$ and $\sigma_2(A_4, A_5, A_3)$; $\sigma_3 = (\sigma_1.A_1, \sigma_1.A_2, \sigma_1.A_3, \sigma_2.A_4, \sigma_2.A_5)$

~~$\sigma_1.A_2 \sigma_2.A_3$~~

DERIVED OPERATORS

"NATURAL JOIN"

$$\mathcal{R}_1 \bowtie \mathcal{R}_2 = \Pi_{A_1, A_2} (\overline{\mathcal{R}_1[A_1 \cap A_2] = \mathcal{R}_2[A_1 \cap A_2]} (\mathcal{R}_1 \times \mathcal{R}_2))$$

	A ₁	A ₂	A ₃
t ₁₁	a ₁₁	a ₂₁	a ₃₁
t ₁₂	a ₁₂	a ₂₂	a ₃₂
t ₁₃	a ₁₃	a ₂₃	a ₃₃
⋮	⋮	⋮	⋮

	A ₂	A ₃	B ₁
t ₂₁	a ₂₁	a ₃₁	b ₁₁
t ₂₂	a ₂₂	a ₃₂	b ₁₂
t ₂₃	a ₂₃	a ₃₃	b ₁₃
⋮	⋮	⋮	⋮

\downarrow
Natural Join Result

A ₁	A ₂	A ₃	A ₄
a ₁₁	a ₂₁	a ₃₁	b ₁₁
a ₁₁	a ₂₁	a ₃₁	b ₁₃
a ₂₂	a ₂₂	a ₃₂	b ₁₂

t₁₁ ⋈ t₂₁

t₁₁ ⋈ t₂₃

t₂₂ ⋈ t₂₂

Example:

"DIVISION"

$\mathcal{R}_1 \div \mathcal{R}_2 \rightarrow$ Take joinable tuples, but take only the attributes in \mathcal{R}_1 but not in \mathcal{R}_2 .

$$\mathcal{R}_1 \div \mathcal{R}_2 = \Pi_{A_1, A_2} (\overline{\mathcal{R}_1[A_1 \cap A_2] = \mathcal{R}_2[A_1 \cap A_2]} (\mathcal{R}_1 \times \mathcal{R}_2))$$

"INTERSECTION" $\mathcal{R}_1 \cap \mathcal{R}_2$

EXTENDED RELATIONAL ALGEBRA

"OUTER JOIN":

$\sigma_1 \bowtie \sigma_2$ (left)

$\sigma_1 \times \sigma_2$ (right)

$\sigma_1 \bowtie \sigma_2$ (full)

- For $\sigma_1 \bowtie \sigma_2$, include $\sigma_1 \bowtie \sigma_2$ and include unjoined tuples in σ_1 and append "Null" to attributes from σ_2 .
- Similarly for $\sigma_1 \times \sigma_2$ and for $\sigma_1 \bowtie \sigma_2$, we do both $\sigma_1 \bowtie \sigma_2$ & $\sigma_1 \bowtie \sigma_2$.

$$\sigma_1 \bowtie \sigma_2 = (\sigma_1 \bowtie \sigma_2) \cup (\sigma_1 \times \sigma_2).$$

"AGGREGATION OPERATORS":

$A_2 \text{g}_{\sigma_1} (A_1)$ { g - count unique
- average
- max
- min
- sum. }

gives the g-aggregation of attribute A_1 grouped according to A_2 .

DATABASE CONSISTENCY

written down in terms of Functional Dependency.

Course-registration

Roll	Name	Cid	Cname

Things which are obvious:

- If roll is same \Rightarrow name has to be same
- If Cid is same \Rightarrow Cname has to be same

We write these as:

$$\text{Roll} \rightarrow \text{Name}$$

$$\text{Cid} \rightarrow \text{Cname}$$

Let $\sigma(R)$ be a table where R is a set of attributes,
 $"\alpha \rightarrow \beta"$, $\alpha, \beta \subseteq R$

means

$$[t_1[\alpha] = t_2[\alpha]] \Rightarrow [t_1[\beta] = t_2[\beta]]$$

$$F: \left\{ \begin{array}{l} \text{roll} \rightarrow \text{name} \\ \text{cid} \rightarrow \text{cname} \end{array} \right\}$$

$$F^+ \left\{ \begin{array}{l} \dots \text{Set of} \\ \text{implied dependency} \dots \\ (\text{e.g. roll, cid} \rightarrow \text{cname}) \end{array} \right\}$$

ARMSTRONG'S AXIOMS:

1. Reflexivity:

$$\alpha \rightarrow \beta \text{ if } \beta \subseteq \alpha$$

2. Augmentation:

$$\text{if } \alpha \rightarrow \beta \text{ then } \alpha \gamma \rightarrow \beta \gamma$$

3. Transitivity:

$$\text{if } \alpha \rightarrow \beta \text{ and } \beta \rightarrow \gamma \text{ then } \alpha \rightarrow \gamma$$

All that
is needed,

Let $\alpha(R)$ be a table where R is a set of attributes,

$$\alpha \rightarrow \beta , \quad \alpha, \beta \subseteq R$$

means

$$[t_1[\alpha] = t_2[\alpha]] \Rightarrow [t_1[\beta] = t_2[\beta]]$$

$F:$

$$\left\{ \begin{array}{l} \text{roll} \rightarrow \text{name} \\ \text{cid} \rightarrow \text{cname} \end{array} \right\}$$

F^+

$$\left\{ \begin{array}{l} \dots \text{Set of} \\ \text{implied dependency} \dots \\ (\text{e.g. roll, cid} \rightarrow \text{cname}) \end{array} \right\}$$

ARMSTRONG'S AXIOMS:

1. Reflexivity:

$$\alpha \rightarrow \beta \text{ if } \beta \subseteq \alpha$$

2. Augmentation:

$$\text{if } \alpha \rightarrow \beta \text{ then } \alpha \gamma \rightarrow \beta \gamma$$

3. Transitivity:

$$\text{if } \alpha \rightarrow \beta \text{ and } \beta \rightarrow \gamma \text{ then } \alpha \rightarrow \gamma$$

All that
is needed.

* Union Rule:

if $\alpha \rightarrow \beta$ and $\alpha \rightarrow \gamma \Rightarrow \alpha \rightarrow \beta\gamma$

Proof:

$$\begin{aligned} & \left[\alpha \rightarrow \beta \right] \text{ and } \left[\alpha \rightarrow \gamma \right] \\ \Rightarrow & \alpha, \alpha \rightarrow \beta, \alpha \rightarrow \gamma \\ \Rightarrow & \alpha \rightarrow \alpha\beta \text{ and } \alpha\beta \rightarrow \beta\gamma \\ \Rightarrow & \alpha \rightarrow \beta\gamma. \end{aligned}$$

{ Suppose $\alpha \rightarrow \beta\gamma$ }

* Pseudotransitivity:

if $\alpha \rightarrow \beta$ and $\beta\gamma \rightarrow \delta$ then $\alpha\gamma \rightarrow \delta$

Proof:

$$\begin{aligned} & \left[\alpha \rightarrow \beta \right] \text{ and } \left[\beta\gamma \rightarrow \delta \right] \\ \Rightarrow & \alpha\gamma \rightarrow \beta\gamma \text{ and } \beta\gamma \rightarrow \delta \\ \Rightarrow & \alpha\gamma \rightarrow \delta \end{aligned}$$

* Decomposition:

if $\alpha \rightarrow \beta\gamma$ then $\alpha \rightarrow \beta$ (and $\alpha \rightarrow \gamma$)

Proof:

$$\begin{aligned} & \alpha \rightarrow \beta\gamma \text{ and } \beta\gamma \rightarrow \beta \\ \Leftrightarrow & \alpha \rightarrow \beta. \end{aligned}$$

- If $A \rightarrow R$ where R is the set of all attributes, then A is a candidate key.

CLOSURE OF AN ATTRIBUTE SET α

α^+ is the set of attributes in R such that $\alpha \rightarrow \alpha^+$ with respect to F^+ , the closure of a functional definition, then α^+ is the closure of the attribute set α .

- If $\alpha^+ = R$, then α is a candidate key.
- If there exists F' (simpler than F) such that

$$F'^+ = F$$

then checking database consistency of F' suffices.
This F' is called a 'canonical cover'.

F' is a canonical cover, if there does not exist F'' such that

$$F'' \subseteq F' \text{ and } F''^+ = F'^+$$

whereas $F'^+ = F^+$.

Example:

$$F = \left\{ \begin{array}{l} A \rightarrow B \\ B \rightarrow C \\ A \rightarrow BC \end{array} \right\} \Rightarrow \text{then } F' = \left\{ \begin{array}{l} A \rightarrow B \\ B \rightarrow C \end{array} \right\}$$

Roll	Name	Cid	Cname
101	AB	CS1	PDS
101	AB	CS2	Algo
102	AG	CS1	PDS

Example : $\text{Roll} \rightarrow \text{Name} \wedge \text{Cid} \rightarrow \text{Cname}$

Extraneous Attribute :

Consider those functional dependencies :

$$FD1 : \alpha \rightarrow \beta$$

$$FD2 : \beta \rightarrow \gamma$$

$$FD3 : \alpha \rightarrow \beta\gamma$$

$FD3$ is implied from $FD1$ & $FD2$, ~~then~~ we can drop γ from $FD3$ and it becomes $FD1$. Here γ is an extraneous attribute.

Obtaining the Canonical Cover F^C of F

1. ~~Apply Union Rule (apparently to get closure)~~
 2. ~~Check for extraneous attributes on L.H.S & R.H.S
remove all the~~
1. Apply Union Rule and add only the unions to F^C
 2. Check for extraneous attributes on LHS & RHS, and remove this.

REDUNDANCIES IN TABLES

Faculty					
	Name	Dept	Salary	Dept-addres	Ph-no
1	AB	CS	1000	Nalanda	123
2	AG	PHY	1000	Takshila	124
3	PM	PHY	1000	Takshila	125
4	NG	CS	1500	Nalanda	126
5	PM	CS	1000	Nalanda	125

Anomalies:

- Update Anomaly:

- Update at one location may cause inconsistency.

- Insertion Anomaly:

- Insertions causing inconsistency.

- Deletion Anomaly:

- Deletion of one row may cause unnecessary loss in information.

* The following decomposition may be helpful.

Faculty

Name	Dept	Salary	Ph.no
(Should be <u>Lossless</u>)			

Department

Dept	Addres

* "Lossy join" — cannot recover from the join

* "Lossless join" — Possible to recover original table after naturally joining the tables.

Lossless Decomposition

If R is split into R_1 & R_2 , then it forms a lossless join, if -

- $R_1 \cap R_2 \neq \emptyset$
- Either $R_1 \cap R_2 \rightarrow R_1$
or $R_1 \cap R_2 \rightarrow R_2$

[which is basically

$$\{R_1 \cap R_2 \rightarrow R_1, R_1 \cap R_2 \rightarrow R_2\} \cap F^+ \neq \emptyset$$

Lossless decomposition makes sense only under a set of functional dependencies.

Dependency Preservation

$R = R_1, R_2, \dots, R_n$ be a split and F be the set of functional dependencies.

$$F_i := \{[\alpha \rightarrow \beta] \in F \mid \alpha, \beta \subseteq R_i\}$$

This decomposition is Dependency Preserving if $(F_1 \cup F_2 \cup F_3 \dots F_n)^+ = F^+$

$$\text{or } (F_1 \cup F_2 \cup F_3 \dots F_n)^c = F^c \dots \dots \dots (1)$$

A Decomposition must be

- Dependency Preserving
- Lossless join.

Algorithm to check if $(F_1 \cup F_2 \cup \dots \cup F_n)^+ = F^+$:

```
for  $\alpha \rightarrow \beta \in F - (F_1 \cup F_2 \cup \dots \cup F_n)$ 
     $\alpha_0 := \alpha$ 
    for i in 1 to n
         $\alpha_i := \text{closure of } \alpha_{i-1} \text{ in } F_i$ 
     $\alpha^+ := \text{closure of } \alpha \text{ in } F$ 
    if  $\beta \notin \alpha_n$ 
        return NOT_DPTM
```

return DP

SCHEMA REFINEMENT

We do a set of decomposition until we reach a "normal form"

Normal Forms:

"In a nutshell, these are non-redundant form"

First Normal Form (1-NF):

A table is in 1-NF if it does not contain any multivalued or composite attribute.

[Just split up into another column if needed]

Boyce Codd Normal Form (BCNF)

A relation $\sigma(R)$ is in BCNF if for every FD $\alpha \rightarrow \beta$ in F^+ the following holds true.

$\alpha \rightarrow \beta$ is trivial or α is a superkey of R .

Example:

Faculty

Name	Dept	Salary	Ph

$\text{Name} \rightarrow \text{Dept}$

$\text{Name} \rightarrow \text{Phone}$

$\text{Name} \rightarrow \text{Salary}$

{ } NOT IN BCNF

Dept

dept	address

$\text{dept} \rightarrow \text{address}$

{ } NOT IN BCNF

Conversion to BCNF:

• Split R into R_1, R_2 for $\alpha \rightarrow \beta$ which are violating, make $R_1(\alpha\beta)$ and $R_2(R - (\beta - \alpha))$
non-trivial
& α is not superkey

Further split R_1 and R_2 if needed for other violating dependencies.

* BCNF decomposition is always a Lossless join but may not be dependency preserving.

* BCNF can still have "transitive redundancy"

Transitive Redundancy

In a set of FD if

$$\alpha \rightarrow \beta \in F^+$$

$$\alpha \rightarrow \gamma \in F^+$$

$$\beta \rightarrow \gamma \notin F^+$$

then F consists transitive redundancy.

Third Normal Form (3NF):

A table $\alpha(R)$ is in 3NF under a set of dependencies F , if $\nexists \alpha \rightarrow \beta \in F^+$,

$\alpha \rightarrow \beta$ is trivial or α is a part of a candidate key.

Example:

Faculty (Name, dept, salary, address, ph) has candidate key (Name, dept), then under

$$F = \{$$

$$\text{Name} \rightarrow \text{Address Phone}$$

$$\text{Dept} \rightarrow \text{Address}$$

}

is not in BCNF but is in 3CNF.

3. CNF decomposition

Take union $\alpha \rightarrow \beta$ and make tables
for each L.H.S.

Now want to find β
so we can find $\alpha \rightarrow \beta$
and $\beta \rightarrow \gamma$ and $\gamma \rightarrow \delta$
so $\alpha \rightarrow \delta$ and $\beta \rightarrow \delta$

Algorithm

such that
(A) \neg variables appear (leftmost) printed?
and next (rightmost) not substituted
 $\} = 7$

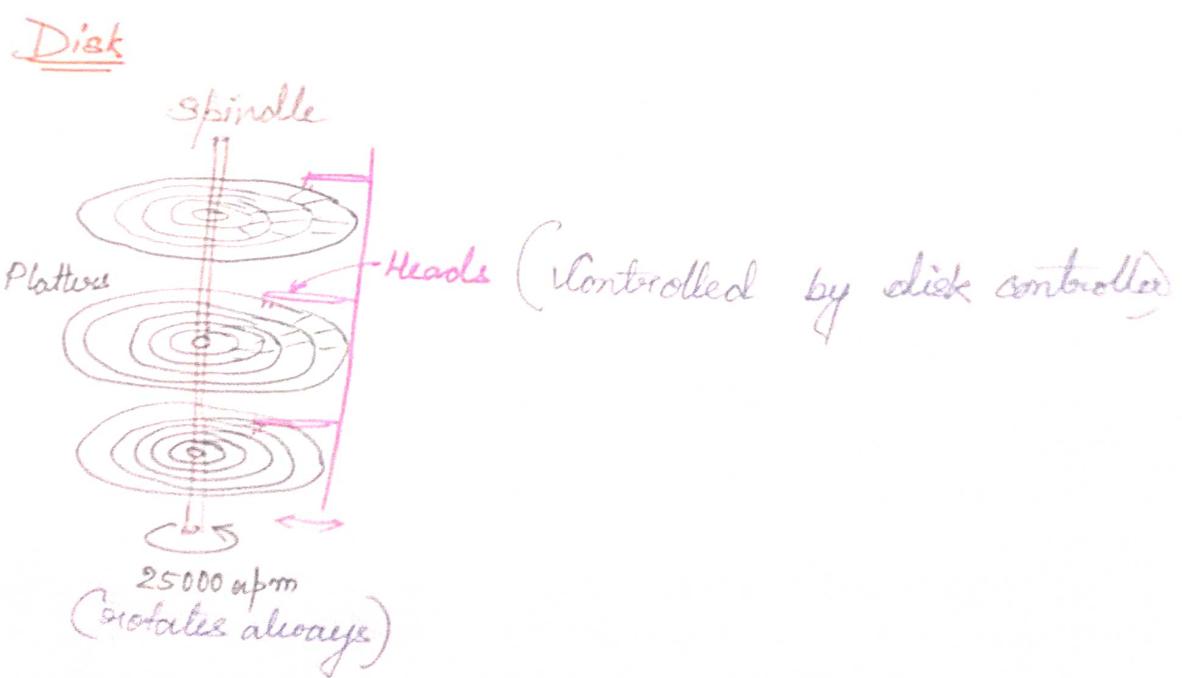
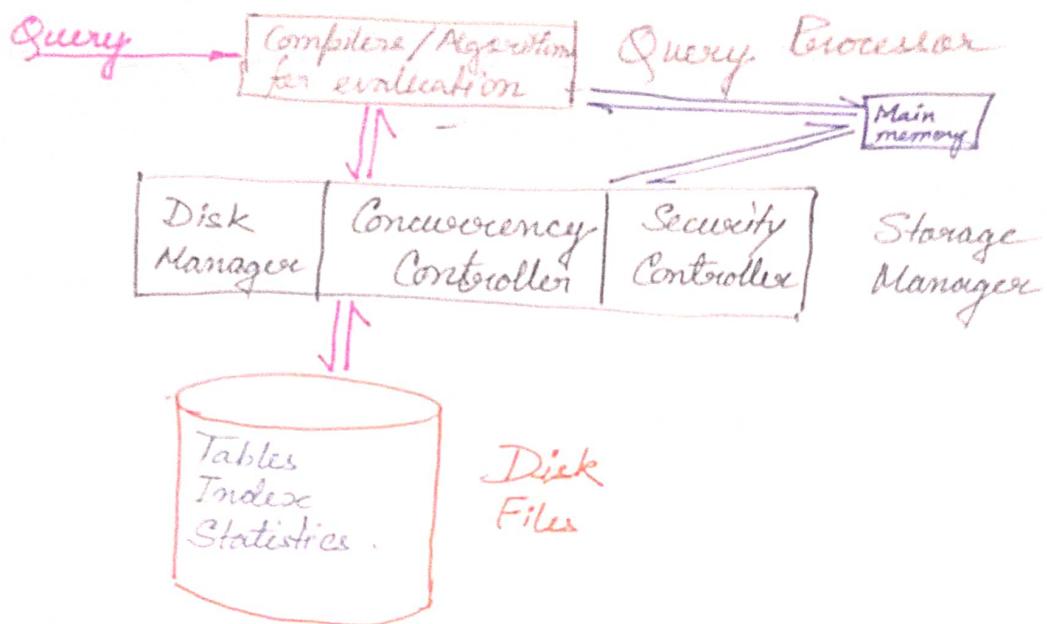
great variable \rightarrow small

variable \rightarrow large

{

BASE is AND TRUE in fact is

PHYSICAL DATABASE DESIGN



Seek-time → Time required for head to move to right track ($\sim 4\text{ ms}$)

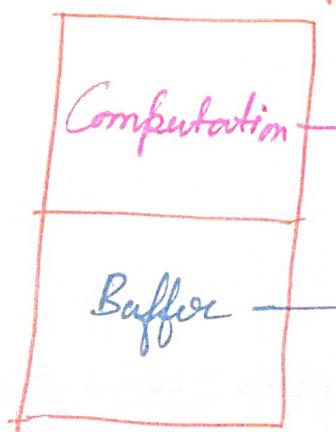
Rotation Latency → Time required for the disk to rotate & the correct sector to arrive ($\sim 2\text{ ms}$)

Disk read/write time is so bad, that any query algorithm is evaluated on the bases of disk read/write.

DBMS (not OS) does the following:

- Disk Controller algorithm
- Memory Management
- Parallel Read/Write.

Main Memory



Computation — CPU is currently using.

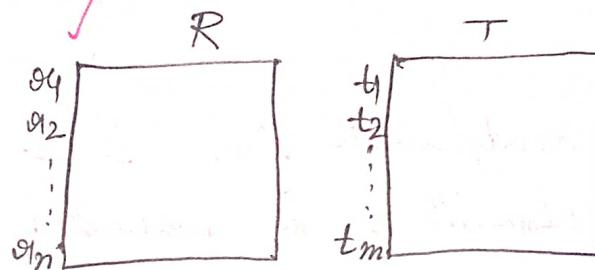
Buffer — Extra data from Disk (like Cache).

- DBMS does not follow LRU scheme because of:

- Consistency

- * Pinned Block

- Optimality



$R \bowtie T$

algo:

for (r from o_1 to o_n)

 for (t from t_1 to t_m)

 if r joins t

 output r, t

LRU is not optimal here.

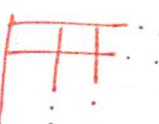
FILE STRUCTURE

Logical Organization

DB



O₁

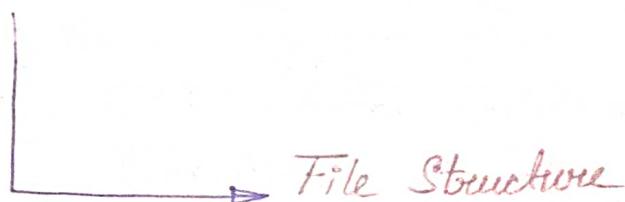


O₂

Hard disk



sectors

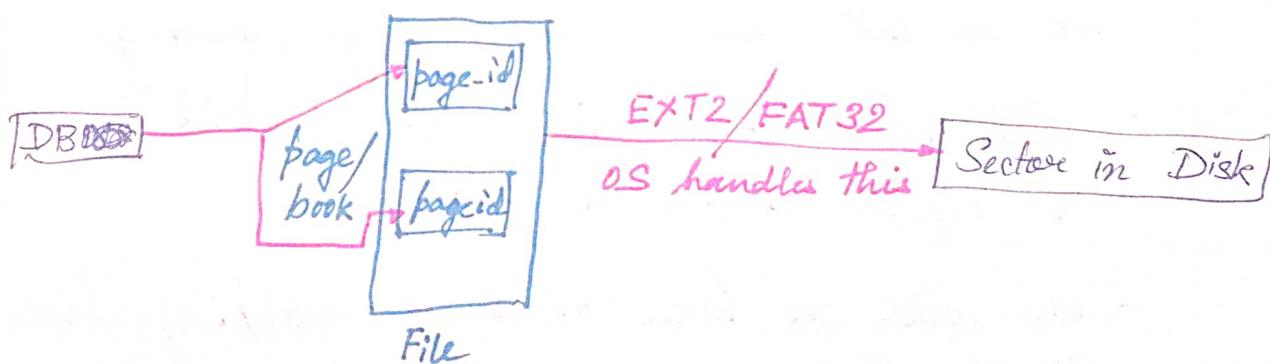


File Structure

- MongoDB has column oriented file structure (efficient for narrow tables)
- We will study tuple oriented file structure

Tuple Oriented File Structure

- Each file has multiple pages/books.



- page id (pid) is mapped to location at disk.
- pid is unique across all files.
- Mapping from DB to ~~file~~ pages is DBMS dependent

Multi Table Clustered File

Files contain multiple tables.

- Good if tables are small. join operations might be easy.

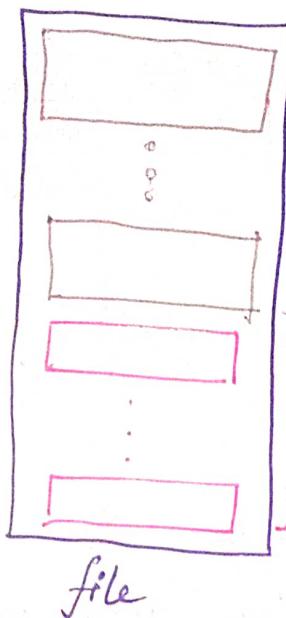
Single table in single file

- disk/Hardwave $\xrightarrow{\text{page}}$ disk — 4 KB (correctness guaranteed by hardware)
- Memory page — 4 KB (no guarantee)
- File page — 512 B to 16 KB (no guarantee)
 - [SQLite page size — 512 B
MySQL page size — 16 KB
Postgres page size — 8 KB]

- Due to the nature of relational algebra, we do not need to follow any ordering.
- Each table identified by a pid and its offset within a pid.
- We will see three kinds of mapping from DB to pages
 - Heap file structure (most popular)
 - Sorted file structure
 - Block file structure.

Heap File Structure

[Assumption: One page file at least one tuple]

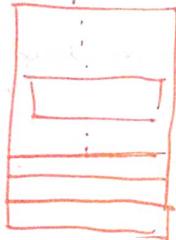


Header (Stores header information)
Pages (of the file)

Tuple Pages (Stores the actual tuple values)

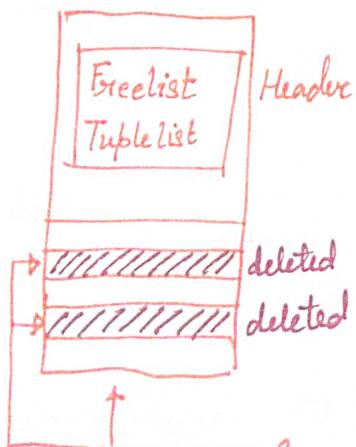
file

Adding Tuples



: ↗ Can add here

Deletion of Tuples



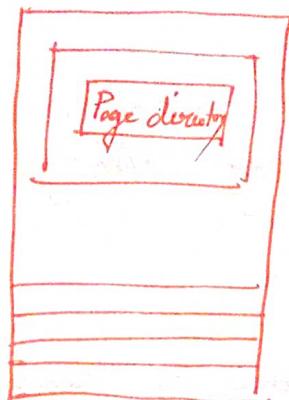
Can add here
(new insertion)

Freelist: Linked list containing the location of all free space
list (including unallocated space)

Tuple list: Linked list containing the location of all occupied space

• Updates are costly, hence this is not done! (o-o, o-o)

- Approach 2: Keep a "Page Directory" in the header

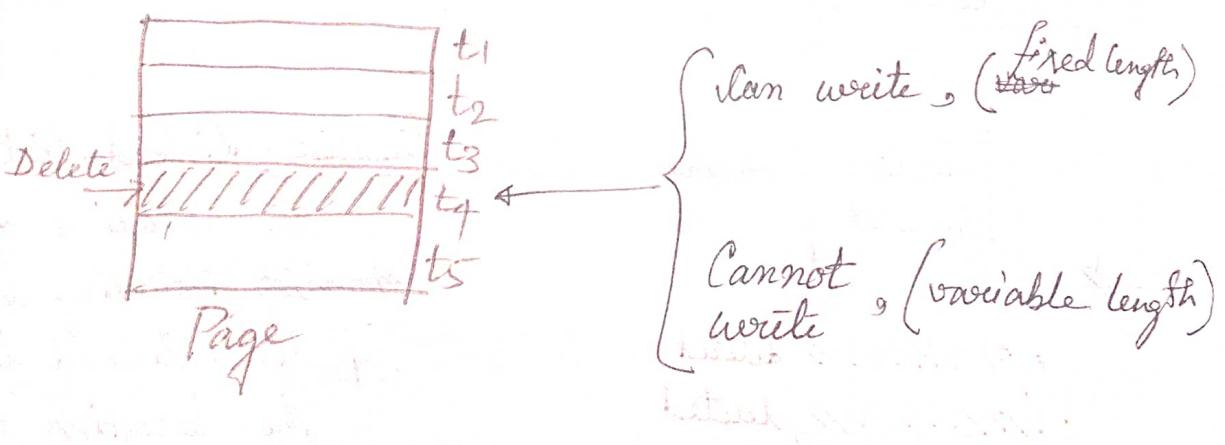


offset	tag
	Free
:	Occupied
	:
	:
	:

Basically
the location
in the table.

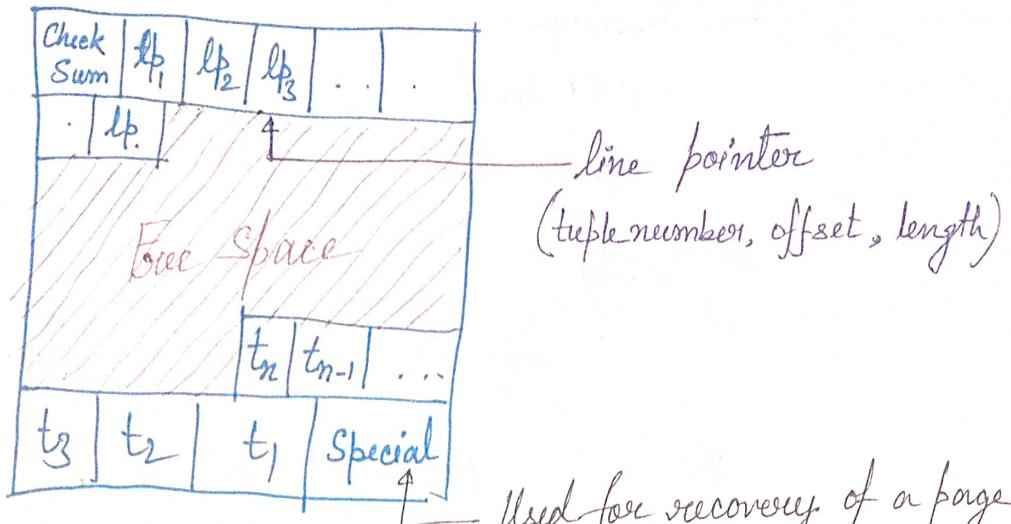
- Turns out that for large files, the memory overhead of Page directory is very less
- This "Page Directory" is small and can be kept in buffer (in memory)
- This, is what is done! 😊

Structure within the pages



- One structure for variable length records is "Slotted Page Structure"

Slotted Page Structure



- Within a page, each tuple is given a number
- When the tuples move around, only the offset changes in the corresponding line pointers

Can try on postgres

> CREATE EXTENSION pageinspect;

> \dget_raw_page('pg_class', 0)

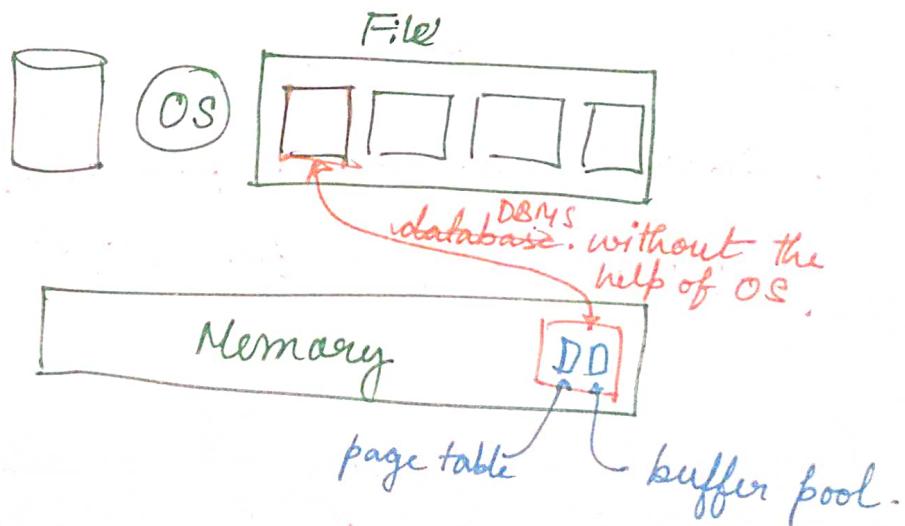
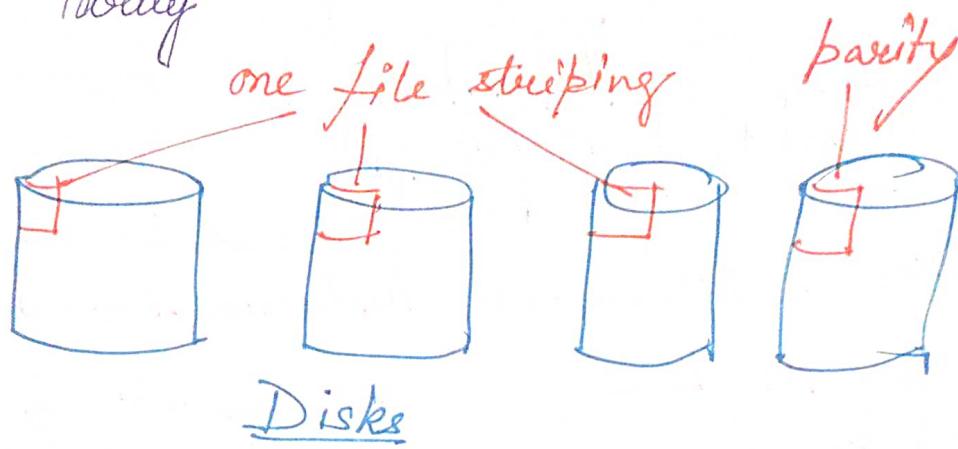
Redundant Array of Independent Disks (RAID)

Advantages of multiple disks

- Parallel Read/Write
- Reliability (by redundancy)

RAID5

- Striping
 - o bit level
 - o page level
- Parity



In postgres you have a command called `CREATE BUFFERPOOL name size`.

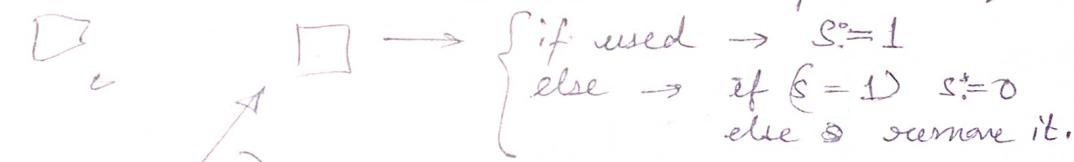
There are things which are needed to be kept in mind (for multithread access or multiple query queries)

- Consistency
- Efficiency

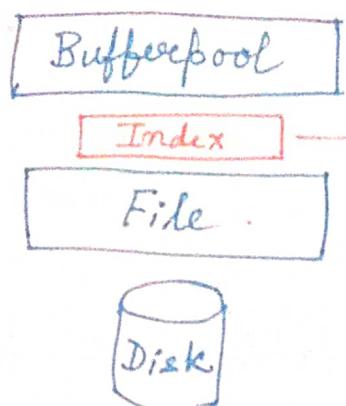
We use something called latch (shorter than locks) which says that you cannot change buffer ~~without~~ without making it durable.

Strategies used in Cache

- LRU (not always optimal)
- Toss immediate (similar to MRU).
- Clock sweep LRU. (approximate version of LRU, easier to implement)



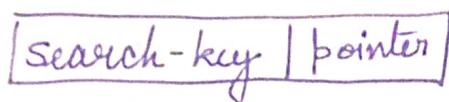
- LRU-K (maintain timestamp of last K accesses, estimate time of next access).



- Stored in disk (as file)
- Anything for fast access.

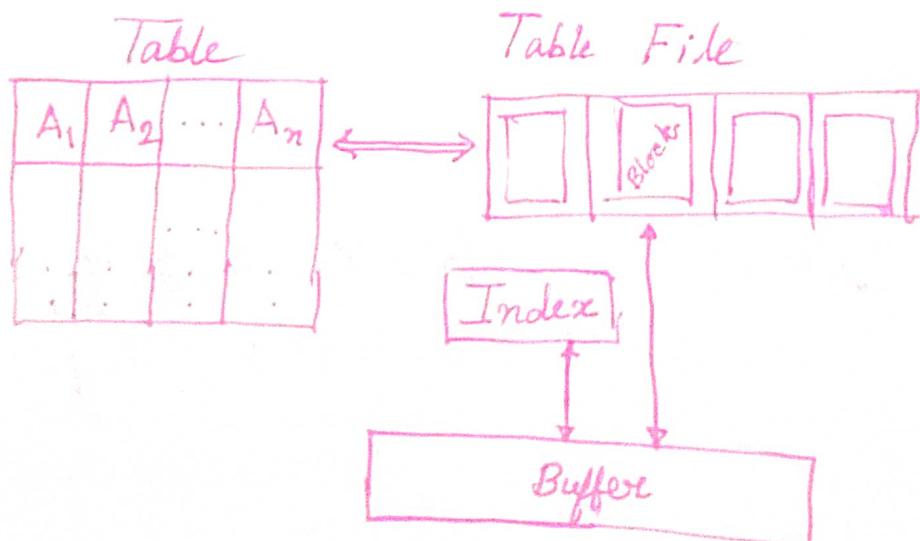
/* OK, sir now moved to slides 😊 */

Index file



- If search-key is sorted → ordered index
- search-keys uniformly distributed into buckets using hash function → hash index

/* More slide stuff */



Structure of Index

Structure of index :

(key, pointer) pairs.

key	pointer
v ₁	p ₁ p ₂
v ₂	p ₃ p ₄
:	

} Stored in files
(Index files)

Command in postgres.

In this episode of 'You have a Command in Postgres, we learn

CREATE INDEX I1 USING AI OF R1

* Note:

- Buffer management of Index files is same as Table files
- Index files needs need to be maintained.
- Index files are usually smaller than table files

Index Evaluation Metrics :

- Access type supported
- Access time.

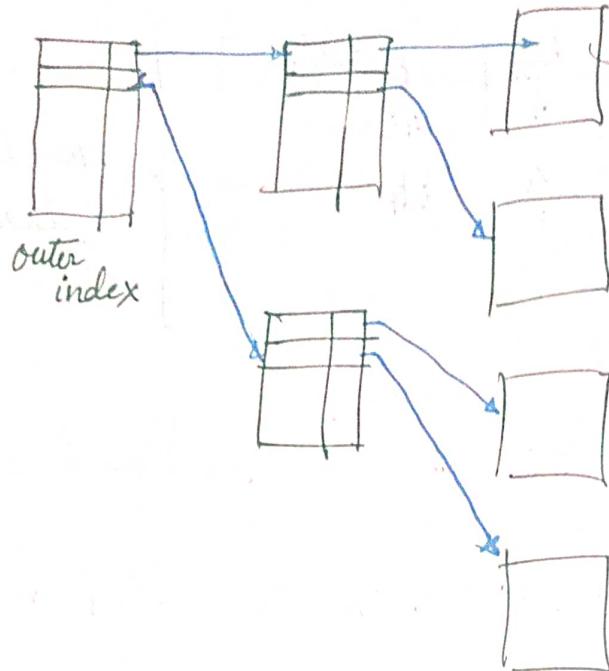
/*

and other things I am too slow too copy 😊

*/

Key insight: What if we keep an index of the index file?

We get a multilevel index!



We get a tree structure (very)

"BST sounds good, but it needs to be stored in memory".

Use B⁺ Trees!

B⁺ tree:

- Spanning more
- Height less
- Each node is the size of a block (disk block)
- Each node is an array of key-pointer pairs.

Properties:

- o Balance: Each path from root to leaf are of the same length (self balancing)
- o Fix n , the order of the tree.
 - * A node has between n & $\lceil \frac{n}{2} \rceil$ pointers
 - * ~~A~~ children.
 - * A node has between $(n-1)$ and $\lceil \frac{(n-1)}{2} \rceil$ values.
 - * n is decided on the basis of disk block size & the key-pointer pair size.

$n \rightarrow$ order, $R \rightarrow$ number of records

Search: ~~$O(\log n)$~~ $O(\log_{\frac{n}{2}} R)$

Delete: $O(\log_{\frac{n}{2}} R)$

Insert: $O(\log_{\frac{n}{2}} R)$

} Number of block I/O
(read/writes)

Values are ~~sus~~.

Observations: (Rather design properties)

"Einstein's children are all less than Einstein"

"Einstein's child are all less than Einstein"

- Leaf nodes point directly to records (or a bucket of pointer).

- All values in a node are sorted.

- pointers pointing to a child from a (pointer, key) pair, all ~~less~~ value keys in the children are less than ~~equal~~ the original key. It is

go values & "to the right" are greater or equal to the original key.

(non-leaf).

"Let us build a tree 😊"

Example: $n=4$, entries = $\{1, 2, \dots, 21\}$
"build bottom up".

* much confusion in class */

Fed up; not writing anything
go read slides! some other source!

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----

Data Access

④ Key $\xrightarrow{\text{fast}}$ Record location in disk.

⑤ Index

Key	P
:	
:	
:	

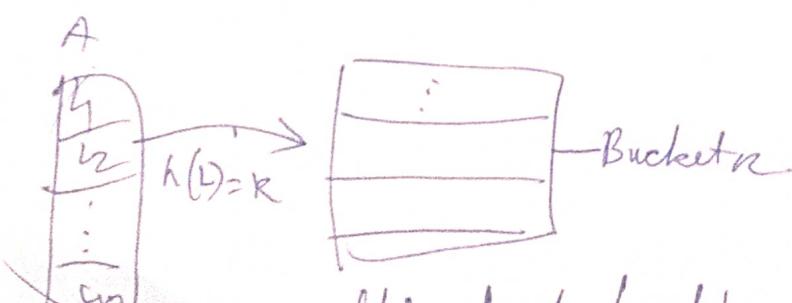
ordered

⑥ Hash Table

(key, value)	Buckets
(key, value)	

(typically size of a disk block)

- The value itself can be a set of pointers (if tuples are small enough) ! Then its called a hash file
- But more commonly, there are pointer pointing to disk blocks containing the data (called hash index).



- (i) apply hash function
- (ii) probe into bucket

Currently accepted best hash function: XXHash3 (given by Facebook)
 Google → FarmHash.

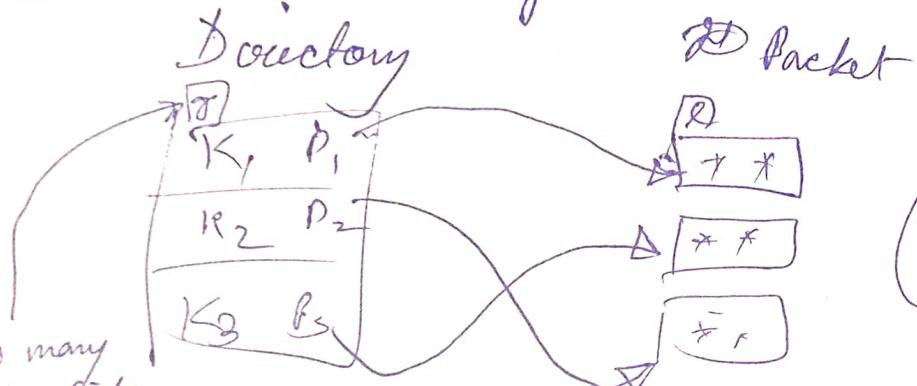
Dynamic Hashing



Overflow Bucket

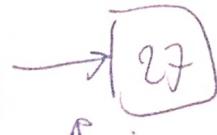
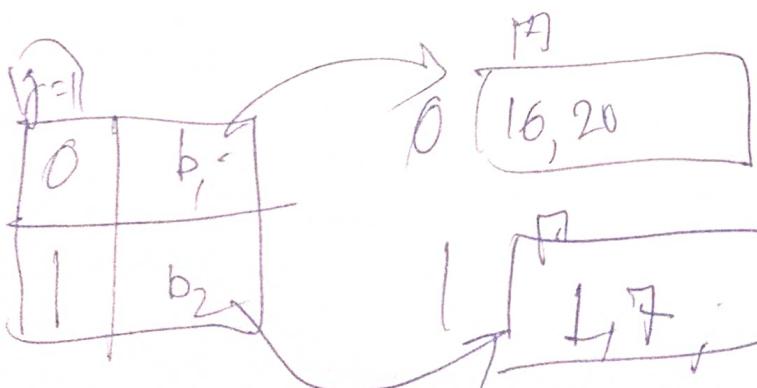
Dynamic Hashing

Extensible Hashing



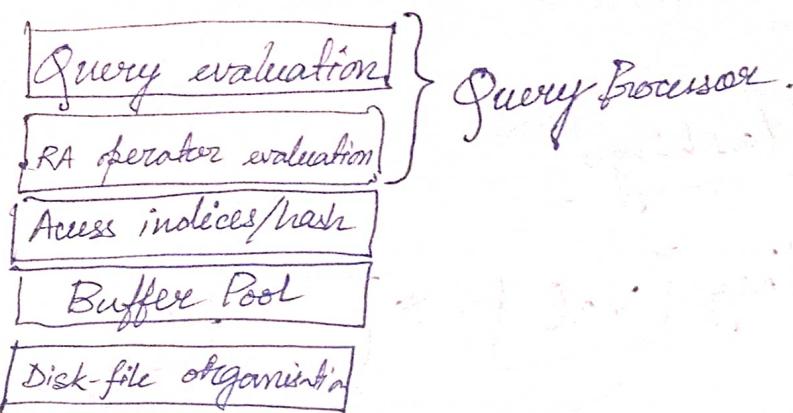
How many
suffix bits
you need to
look at

(can split if
 $l < g$)



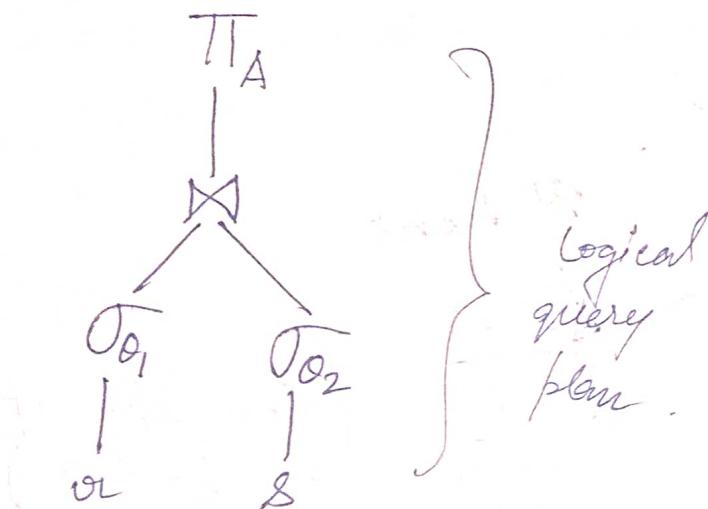
when overflow occurs,
double g.

DBMS



Parse tree:

$$\Pi_A(\sigma_1(r) \bowtie \sigma_{O_2}(S))$$



Cost: Number of disk I/O's.

Optimizer tries to find the least cost tree (NP-Hard)

Physical Execution Plan

(Annotated parse tree containing algorithms & indices.)

Disk Cost

Disk Cost

- (# Seek) * (avg seek time)
- (# Blocks read) * (avg read time)
- (# Blocks written) * (avg. block write time).

Evaluation of operators

O₀

O → match selection

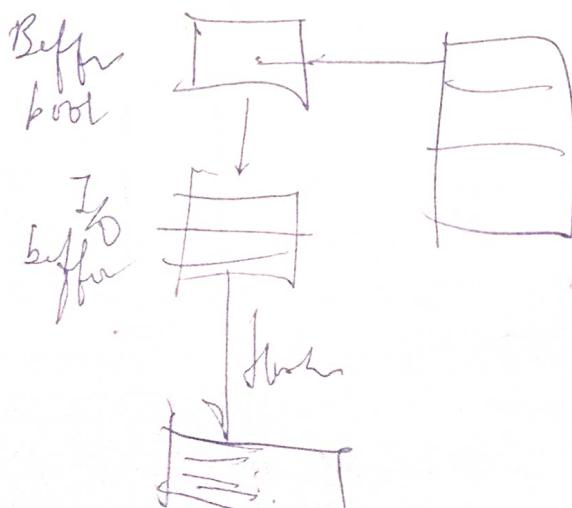
$$(A = 10)$$

O → Range selection

$$A > 10$$

O → complicated selection $\otimes Q_1 Q_2 \dots Q_n$

Match selection



Range Range query
(easy using $B+$ tree)

Conjunction -

$$\sigma_{q_1 \wedge q_2 \wedge \dots \wedge q_n}(\pi)$$

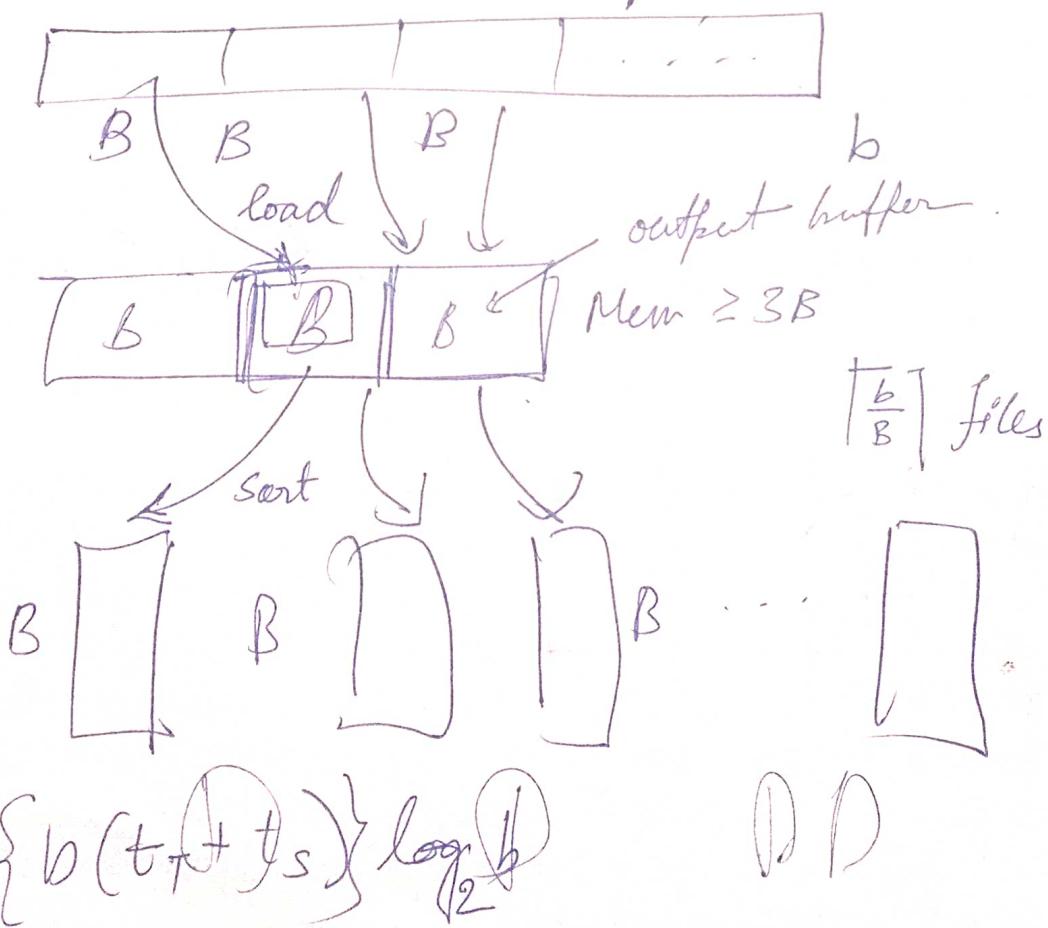
check for the strictest Q_i 's, take only satisfied ~~parts~~ ^{tables} to buffer, and see other conditions in buffer.

Sorting

Memory — Megabytes

Files — Terabytes

Hack: Chop up the memory into pieces.



Join Operation

④ Stupid nested loop join
— as name implies.

⑤ Block Nested Loop joint.
— just keep loop on blocks.
— indexed file in inner loop.

⑥ Merge - join
— merge on sorted attributes
(if indexed, you move down the index,
instead of the attribute value).
~~index~~

⑦ Hash - join :

[Sorts cannot be in place, or anything.
→ other queries may be reading the same
file]

DB MS
20/08/23

Query SQL

Query Rewrite RA

Query Transformation RA

Query Plan

Computationally
hard.

Optimization
Algorithm

Cost Estimate

Statistically
hard

Best plan -
(NP Complete)

Prefer $\sigma_Q(\sigma_{O_2}(E))$ rather than $\sigma_{O_1 \wedge O_2}(E)$

Query Optimization

1. Generate Equivalent Algebra Expressions.
2. Annotate with execution plan
3. Select the best plan

3.1 Use some Heuristic Rules

- move down σ/π
- sequentialize, conjunctive predicates

"Selective of a predicate for a table is

select (output table)
(original table)

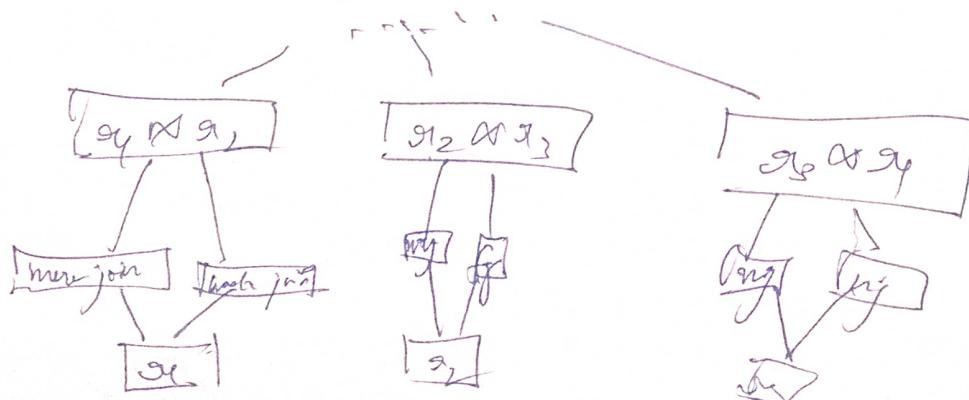
- replace $\sigma-X$ by $\exists X$

3.2 Cost based optimization -

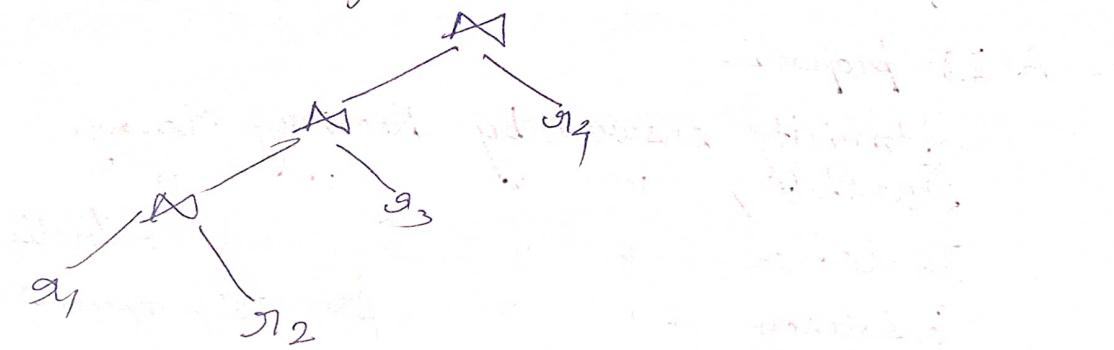
(Combination of search & heuristic.
Join operations are explored, σ & π
(\exists) are not explored).

Selengor style optimization (bottom up) :

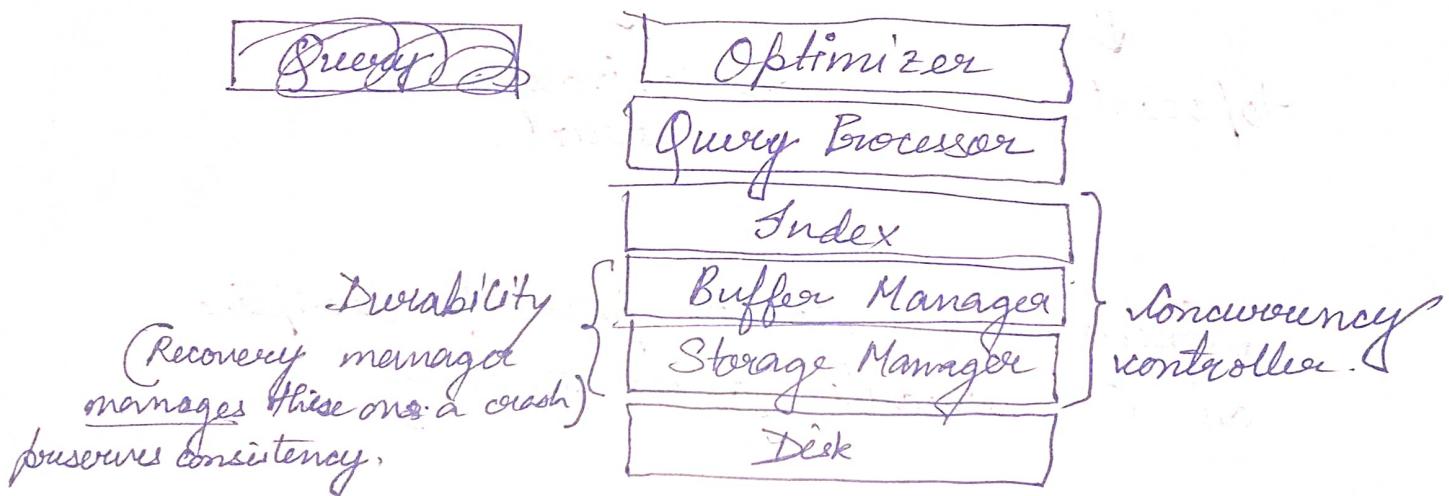
$\sigma_1 \wedge \sigma_2 \wedge \sigma_3$



Produces left-diff. join tables



Interesting join orders — Some orders already have something sorted



Transaction Concept

- Set of instructions that must be executed as a unit
- ACID properties:

Atomicity	ensured by	Recovery Manager
Durability	" "	" "
Isolation	" "	Concurrency Control
Consistency	" "	Forward Progress

Precedence graph = Conflict
is dag serializable

toposort → equivalent serial.

Concurrency Controller puts up discipline
on OS schedule (called protocols)

TRANSACTIONS

Transfer 1000 Rs from account A to B

Transaction

START
 Read(A)
 $A = A - 1000$
 Read(B)
 $B = B + 1000$
 Write(A)
 Write(B)
 COMMIT

- Transactions start with START
- Transactions end with COMMIT / ABORT

Schedule
Conflict equivalence } Much
 Concepts

Dependency Graph



T_1 precedes T_2 in a conflicting pair of operations.

"Pessimistic protocols are popular & easy to implement"



Popular : Locking Protocol

Locking Protocol

lock - $X(A)$ - unlock
 lock - $S(A)$ - unlock
 shared.

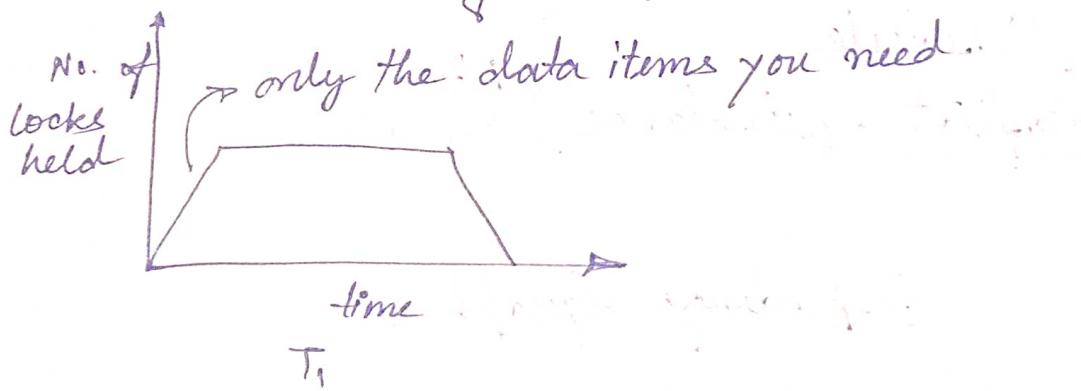
T_1	A	X
S	\checkmark	x
X	x	x

"There is a Lock manager" which maintains a Lock table"

maintains some kind of priority queue.

Transactions

Two Phase Locking (2PL)



Read(A)

Read(B)

$A = A + 1000$

$B = B - 1000$

Write(A)

Write(B)

lock - $X(A)$

lock - $X(B)$

Read(A)

Read(B)

$A = A + 1000$

$B = B - 1000$

Write(A)

Write(B)

Unlock(B)

Unlock(RA)

"Any protocol which is Two-phase locking are conflict serializable".

→ may lead to cascaded rollback & deadlocks

→ serial → sorted according to unlocking time.

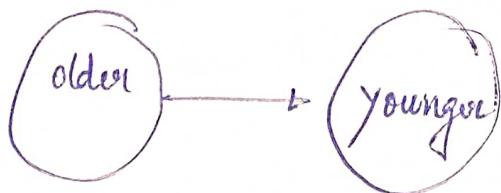
Timestamp Based Concurrency Control

time-stamp order = serializability order.

(Read)

Multi timestamp (α) \rightarrow Time stamp of youngest transaction that has written it successfully.

[There can never be cycles in the precedence graph].



CONCURRENCY CONTROL.

SQL
Normalization } Concepts .

Disk Management .

Buffer Management

Disk Controller .

File Organization .

Indexing → Def / B⁺ tree
(including dynamic hashing)
↓ extensible

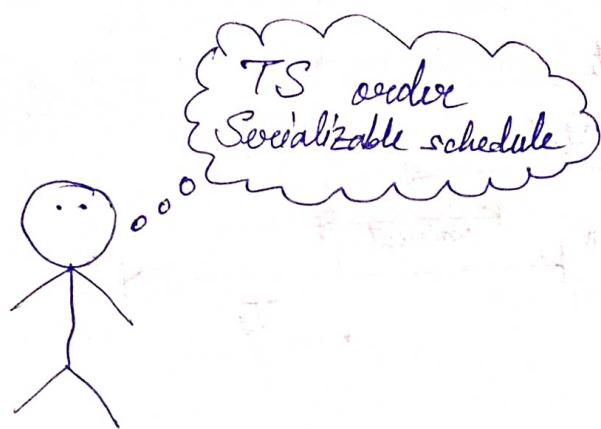
heavy

Query Processing → { Join ,
select
mergesort . }

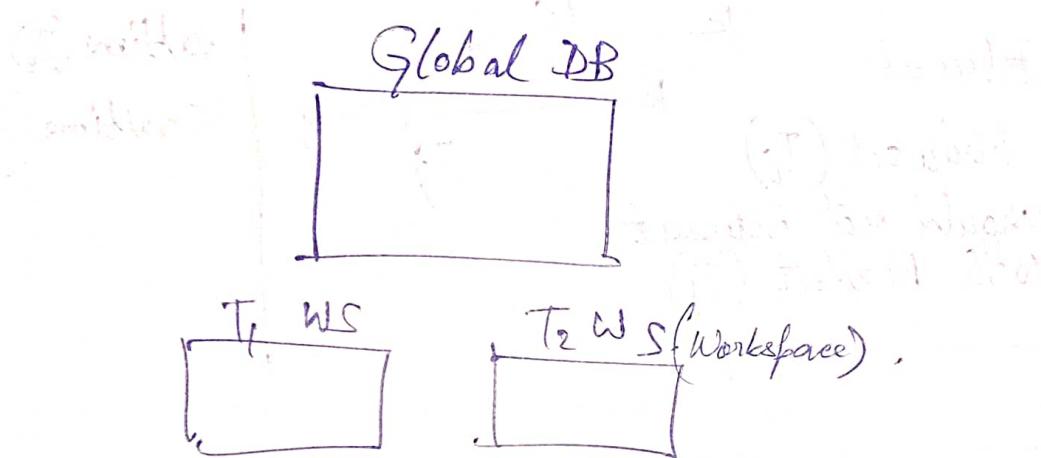
Query Optimization } - Basic notions of
cost based optimization
- Parse tree
- which tree is better
which tree has better
- to how calculate costs

heavy

Transaction { ACID Properties
Concurrency
Serializability, Recoverable,
Cascaded roll back in schedule .
→ 2PL
→ Timestamp
→ Validation
Recovery → log based ,



Concurrency
Controller

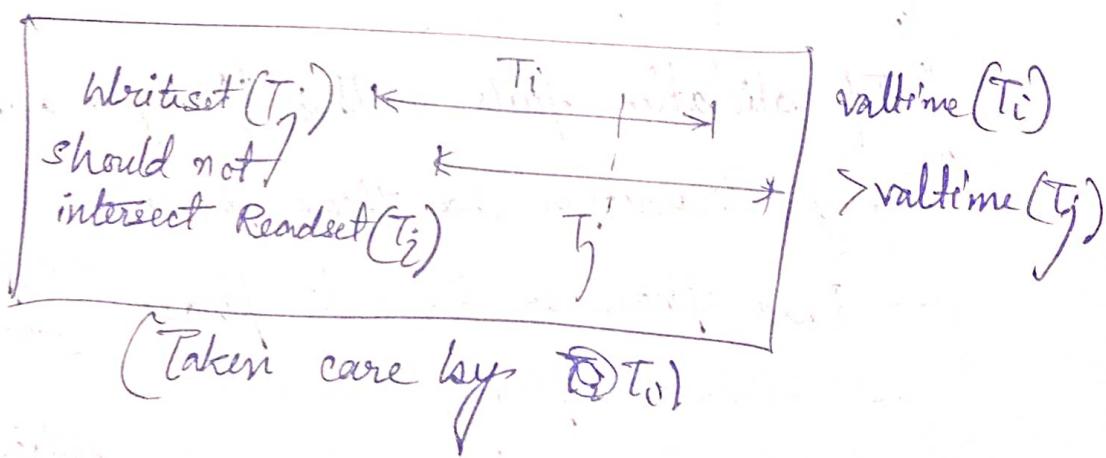
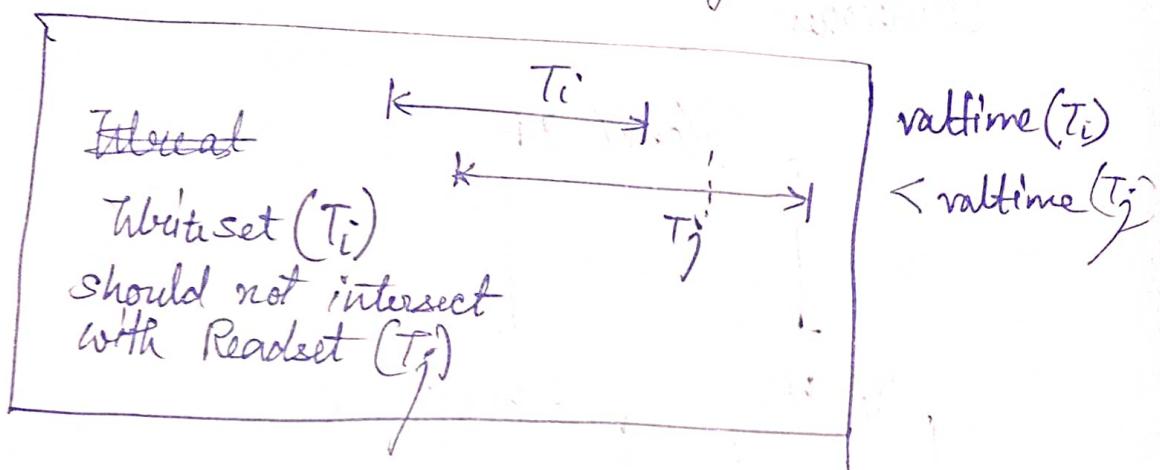
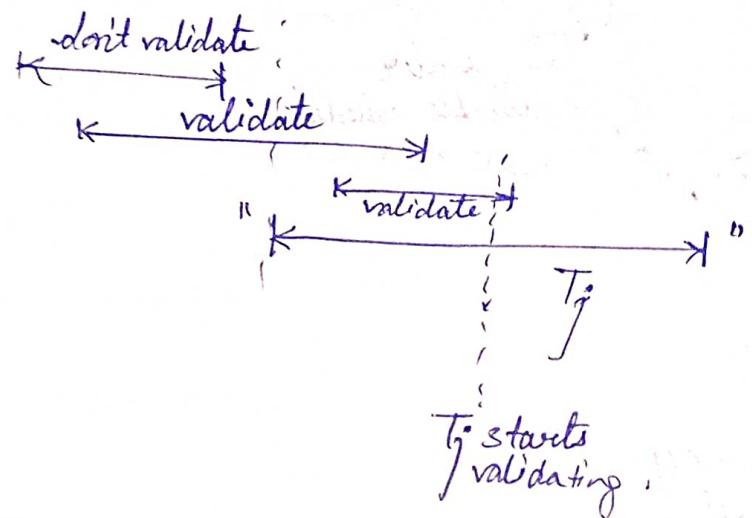


"If validation fails, rollback (i.e. give a new timestamp)"

"Every transaction has three phases".

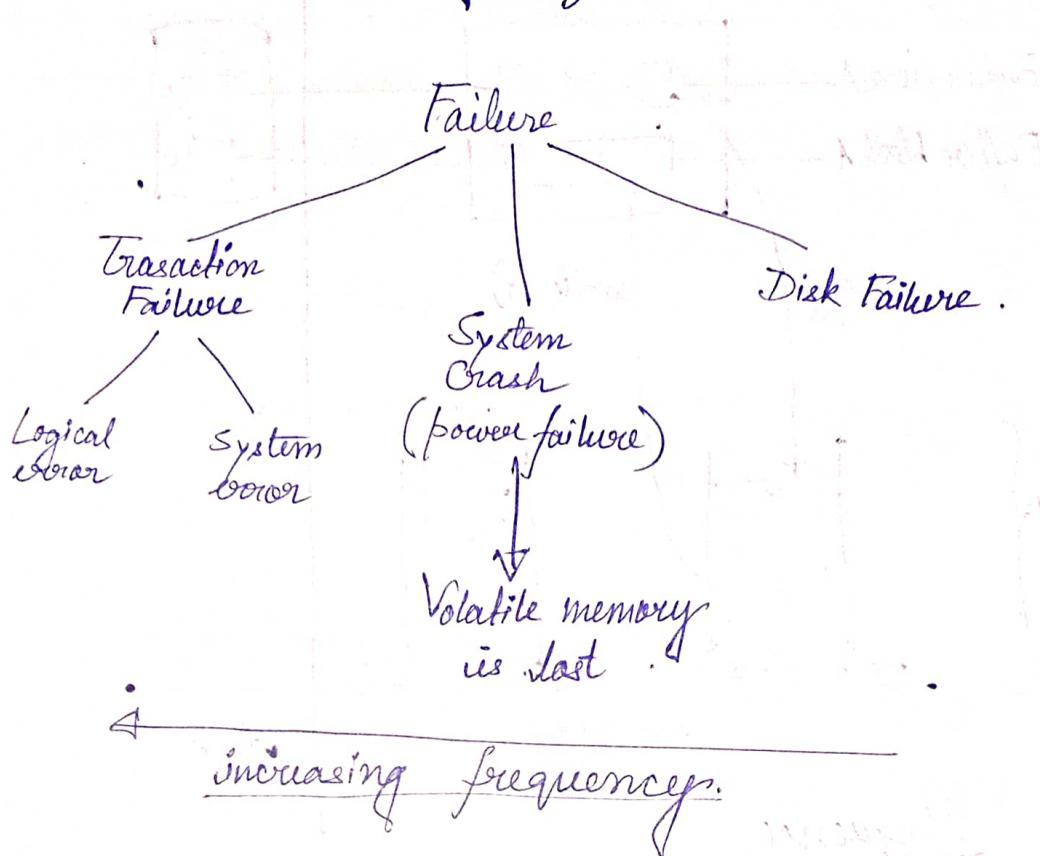
- Read phase & Execute phase (Copy to local workspace)
- Validation phase (Atomic) } entirely atomic
- Write phase (Atomic) } sometimes

"Now validation timestamp is THE timestamp instead of start timestamp. Thus, longer transactions no longer starve, at the cost of the extra workspace!"



"Snapshot → any transaction trying to write to a data item, creates a version."

Recovery System

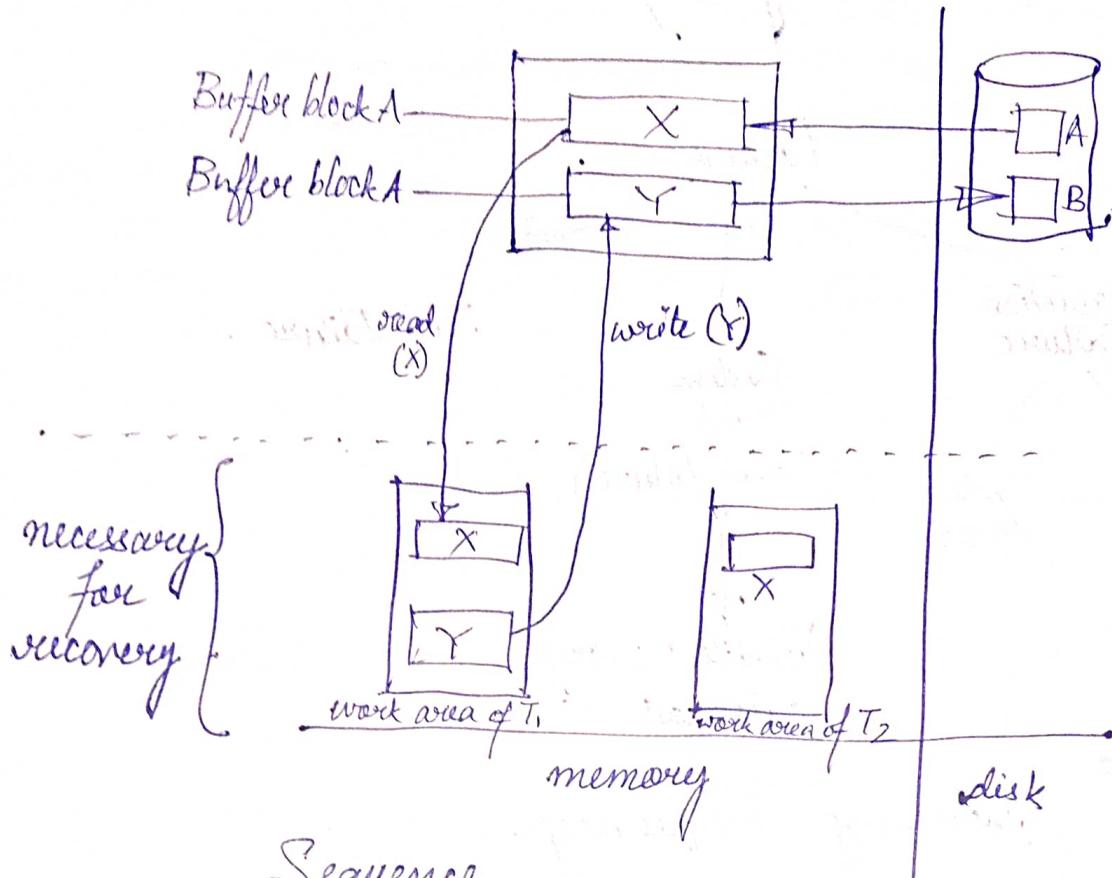


Recovery System : (Atomicity & Durability)

Recovery Algorithm → What to do in normalcy
What to do on crash.

Stable Storage: → Without change, data stays here till infinity, without failure.

→ Implemented by multiple disks, possibly separated geographically.



Sequence

- First perform input
- Then read/write
- Then output .

- Log based recovery.**
- Before writing, write to log file about what I want to write.
 - The log file must be written to stable storage, and only then is it allowed to actually perform the write.

Log Records

- $\langle T_i, \text{start} \rangle$: Transaction T_i starts.
- $\langle T_i, X, V_1, V_2 \rangle$: Transaction T_i wants to write X from old value V_1 to new value V_2 .
- $\langle T_i, \text{commit} \rangle$: T_i finishes the last statement.

Changes: Local Workspace \rightarrow Global Workspace \rightarrow Disk
Two approaches

- Immediate database modification
- Deferred database modification.

New definition of Commit: All previous log records have been written.

Idempotent {
undo(T_i):
 → change value to old value
 → add compensatory record that undo has been done

redo(T_i):
 → keep to same new value
 → compensatory record of redo

undo(T_i) \rightarrow undo(L_{i_1}), undo(L_{i_2}) ... undo(L_{i_k})

redo(T_i) \rightarrow redo(L_{i_0}), redo(L_{i_1}), ... redo($L_{i_{k-1}}$)

When we want to rollback a transaction.

When crash occurs

- Look at the log records that were written just before crash
- undo transactions, having a start but no end commit.
- redo transactions having both a start & commit.

The compensatory log records are used when system crashes while recovering.

Checkpoints

- Hold all inactive transaction
- Output all log records from main memory to stable storage
- Output all modified buffers to the disk
- Write a log record (checkpoint L) onto stable storage where L is a list of all inactive transaction
- All updates are stopped while doing checkpointing.