

## ① \* Concurrency Control in Distributed Databases



### ② Locking Protocols - lock manager deal with replicated Data

#### ③ Single Lock Manager

- One site "S<sub>i</sub>" maintain locks.

Transaction needs to lock data item, first send request to "S<sub>i</sub>" and wait until grant request.

बिपीनकुमार अम. पटेल

पदस्थ - अधिकारीजनों के लिये होरिं प्रबंधक वर्षन

Bipinkumar M. Patel  
Member of the Zonal Manager's Club for Agents  
828 83A

Transaction need data from any sites but write perform on all sites

- Adv i.) Simple Implementation - one message for lock req.
- ii) Simple deadlock handling

DisAdv i.) Bottleneck - "S<sub>i</sub>" become a bottleneck  
ii) Vulnerability - Concurrency controller is lost if "S<sub>i</sub>" fails.

- b) Distributed lock Manager - lock manager function is distributed over several sites.
- each site maintain lock manager and resides at site S<sub>i</sub>.
- when transaction wishes to lock data item φ that is not replicated

→ It is simple in implementation and reduce the degree of bottleneck.

→ However, deadlock handling is much more complex.

#### c) Primary Copy.

- Copy of φ must be reside in precisely one site
- when Transaction needs to lock a data item φ, it request a lock at site φ. The primary copy enables concurrency control for replicated data. If φ fails, φ is inaccessible even though it is present on other sites.

#### d) Majority Protocol

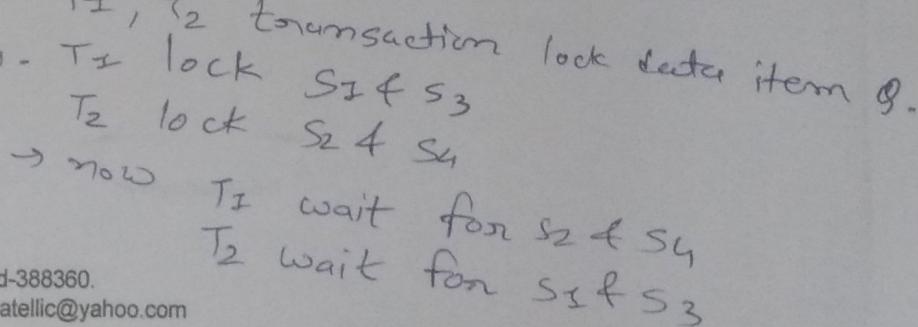
- φ is replicated in 'n' different sites, then lock request message must be sent to more than one-half of the n site in which φ is stored.
- The transaction does not operate on φ until it has successfully obtained a lock on a majority of the replicas of φ.
- Adv - Easy to deal with site failures.

DisAdv i.) Implementation - at least  $2(\frac{n}{2} + 1)$  message for lock request  
ii) Deadlock handling - suppose φ is available on S<sub>1</sub>, S<sub>2</sub>, S<sub>3</sub>, S<sub>4</sub>. T<sub>1</sub>, T<sub>2</sub> transaction lock data item φ.

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणंद-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.



Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

## ④ Biased Protocol -

- shared locks are given more favorable treatment than requests for exclusive locks.
- i) shared locks - request a lock on  $Q$  at one site only
- ii) exclusive locks - request a lock on  $Q$ , at every site  
→ less overhead on "read". Working good when "read" is more frequent than "write".

## ⑤ Quorum Consensus Protocol

- generalization of majority Protocol.
- It assigns read and write operations on an item to two integers called read Quorum " $\Phi_r$ " and write Quorum " $\Phi_w$ ".  
$$\boxed{\Phi_r + \Phi_w > s \text{ and } 2 * \Phi_w > s}$$
- To read and write enough replica must be locked.
- Small reads need to obtain fewer locks but write will be higher.
- This protocol can simulate the majority Protocol & biased Protocols

## ⑥ Time Stamping

- Unique timestamp is given to each transaction
- Two methods for unique timestamp generation
  - ① Centralized - single site distributes timestamp (logical counter used)
  - ② Distributed - each site generate local timestamp using either logical counter or the local clock.
- Problem - Some site generate timestamp fast  
sol<sup>n</sup> - define within each site  $s_i$  a logical clock ( $LC_i$ ) which generate unique local timestamp. Logical clock can be implemented as a counter that is incremented after a new local timestamp is generated
- $s_i$  advance its logical clock whenever a transaction  $T_i$  with timestamp  $\langle x, y \rangle$  visits that site and  $x$  is greater than the current value of  $LC_i$ . for [Synchronization].

198

### ③ Replication with Weak Degrees of Consistency

- **[Master-slave]** replication - update at a 'primary' site and automatically propagates update to other sites.
- Some transactions do not obtain locks at remote sites
  - Replica reflect all update of transactions up to some transaction to the primary data
  - This transactions are in serialization order.



बिपीनकुमार अम. पटेल

www.achivedit.com द्वारा संस्कृत कर्म

Bipinkumar M. Patel

Member of the Zonal Manager's Club for Agents

822883A

- Database configured to propagate update immediately after they occur at the primary or updates only periodically.  
This replication use in central office to branch offices of an organization.

- Another use, creating a copy of the database to run large queries so that queries do not interfere with transactions. Updates should be propagated periodically - every night.

**Multimaster** replication updates are permitted at any replica of data item and automatically propagated to all replicas.

Transaction update local copy & system update's other replicas.  
One way of updating replicas is to apply immediate update with two-phase commit, using one of the lock manager generally biased protocol used.

Many database systems provide an alternative form of update. They update at one site, with lazy propagation. Lazy propagation allow transaction processing to proceed even if a site is disconnected from the network.

④ Update Primary site and it will update other replica lazily.

- Problem some site have old values, serializability problem can be occurs

⑤ Update only replica and propagated other sites.

Problem since the same data item may be updated concurrently at multiple sites.

## ④ Deadlock handling

- Deadlock prevention may result in unnecessary waiting & rollback.

→ Each sites have local wait-for graph

- In deadlock detection, the main problem is deciding how to maintain the wait-for graph.  
SOM

→ These local wait-for graphs are constructed in the usual manner for local transactions and data items. When a transaction  $T_i$  on site  $S_1$  needs a resource in site  $S_2$ . If resource is held by transaction  $T_j$ , the system insert edge from  $T_i \rightarrow T_j$  in the local wait-for graph of site  $S_2$ .

→ If local wait-for graph has cycle which means deadlock but no cycle which not mean no deadlock.

→ Some time Union of two sites wait-for graph has cycle.

\* Centralized deadlock detection, system constructs and maintain global wait-for graph in a single site: deadlock detection

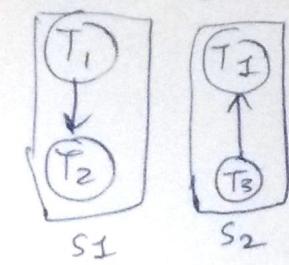
- Some time controller find deadlock but it is not that is because of "false cycles".

The global wait-for graph is constructed if  
i) new edge is inserted or removed  
ii) periodically  
iii) Coordinator needs to invoke cycle-detection Algo.

i) False Cycle - occur due to message delays.

ii) Actual Deadlock - find victim and coordinator notify all the sites and roll back that transaction.

falsecycle



→ If  $T_2$  request  $T_3$  in site  $S_2$  and  $T_1 \rightarrow T_2$  is complete

→ But before delete edge between  $T_2$  &  $T_2$ , system add edge between  $T_2$  and  $T_3$

→ Coordinator find cycle  $T_1 \rightarrow T_2 \rightarrow T_3 \rightarrow T_2$   
Actually deadlock is not present

## \* Availability



- functioning even during failure referred to as robustness.
- All failures handle with different ways
- message loss → Retransmission
- link failure - alternative route
- If system not find other route means network partition
- System find failure but can not find which type of failure

बिपिनकुमार अम. पटेल

सदस्य - अधिकारी के लिये लेत्रिय प्रबंधक कल्पना

Bipinkumar M. Patel

Member of the Zonal Manager's Club for Agents

828 83A

\* Suppose that site S1 has discovered failure

- ① Abort such transaction, locks on data at sites that are still active.
- ② wait until the site to become accessible again
- ③ when data objects are replicated then read & updates
- ④ when failed site recovers, update all changes.
- ⑤ If Replicated data are stored at failed site, catalog should be updated, so that queries do not reference the copy at the failed site when site rejoins, care must be taken to ensure data consistent.
- ⑥ If failed site is a central server then election must be held to determine new server.

→ In particular, these situations must be avoided to ensure consistency.

- ① Two or more central servers are elected in distinct partitions
- ② More than one partition updates a replicated data item.

## \* Replication Protocol design

### ① Majority-Based Approach

- Each data object store with it a version number to detect latest value
- lock request is same as majority based lock

② Read operation look at all replicas and read the value that has the highest version number.

Write operation look at all replicas and find highest version number and increment one and store into all replicas. which has obtained locks, and set the new version number.

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणंद-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.

- If above two rules are violated, the transaction must be aborted.

Here two phase commit protocol can be used, on the sites which are available.

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

- Here reintegration is not required because many replicas are updated so read can be identify latest value easily.

## ② Read One, Write all

→ Special case of quorum consensus -  $\text{readquorum} = 1$   
 $\text{writequorum} = n$

→ In this case, if single site containing a data item fails, no write to the item can proceed since the writequorum not available.

→ Any available replica can be read and a read lock is obtained. A write operation is slipped to all replicas if a site is down, the transaction manager proceeds without waiting for site recover.

→ Reintegration action is required for down sites.  
→ This will work fine for no network partitioning system.

## ③ Site Reintegration

- Use for repair sites
- When failed site recovers, it must reflect changes.
- It is complicated
- entire system halt temporarily while failed site rejoins
- There are techniques have been developed to allow failed sites to reintegrate and concurrent updates to data items.
- Before read or write lock is granted system ensure that it has caught up on all updates to the data item.

## ④ Comparison with Remote Backup

- actions such as concurrency control and recovery are performed at the ~~single~~ site, only data and log records are replicated at the other site.
- Remote Backup avoid Two phase commit and its overheads.
- It provides high availability at lower cost

## ⑤ Coordinator Selection

- Some time coordinator fails and execution stops. One way to continue execution is by maintaining a backup to the coordinator, which is ready to assume responsibility if the coordinator fails.
- Backup coordinator has all information which is required in execution but it can not ~~not~~ take action that affects sites.
- If coordinator fail, backup coordinator take charge. Prime advantage is continue processing without any interruption.
- If backup ~~is~~ is not present, then newly appointed coordinator would have to seek information from all sites.
- If it get only few information then it may be abort several active transactions, and to restart them under the control of the new coordinator.

- backup coordinator avoid substantial amount of delay in new coordinator schema.
- Disadvantage is the overhead of duplication of tasks. Furthermore Coordinator and backup need to communicate regularly for synchronization.



भारतीय जीवन बीमा निगम  
LIFE INSURANCE CORPORATION OF INDIA

**बिपीनकुमार अम. पटेल** - In the absence of backup coordinator, election algorithms enable new coordinator. This algo. require unique identification number be associated with each active site in the system.

**Bipinkumar M. Patel** Member of the Zonal Manager's Club for Agents 828 83A

Bulky algo. for election

- identification number of site "Si" is "i" and site is active with the largest identification number. Whenever coordinator fails, the algo. choose active site which have largest identification number.
- In this case, si sends request to coordinator within a prepecified time interval T coordinator should response if not si will take charge and send election message to every active site that has a higher identification number. If sites not respond, for specific time T then it will elect itself. and inform all sites ~~that of~~ ~~other~~ ~~coordinator~~ with identification number less than i. If si does receive an ~~no~~ answer, it begin a time interval T', to receive a message informing it that a site with a higher identification number has been elected.
- If si receives no message within T', then it assumes the site with a higher number has failed and restart algo.
- After a failed site recovers, it immediately forces all sites to change coordinator if no higher number present. That's why it's called bulky algo.

## ⑥ Trading off Consistency for Availability.

- If a network failure results in more than two partitions, no partition may have a majority of sites. The wait-all protocol provide availability but not consistency.

**CAP theorem** → Distributed Database can have at most two of the following properties ① Consistency ② Availability ③ Partition tolerance

→ Consistent Data if result is the same as if the operation were executed on single site.

कार्यालय : परबड़ी के पास, पो. समरखा.  
निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आर्णद-388360.  
फोन : (नि) (02692) 256673 (मो) 98242 32755.

→ In any large scale distributed system partition can not be prevented and as a result either availability or consistency has to be sacrificed.

- Bank system require ACID properties
- Social-networking system require BASE properties:
  - Basically available
  - soft state
  - Eventually consistent
- Primary requirement is availability, updates should be allowed, even in the event of partitioning. Soft state refers to the property that the state of database may not be precisely defined, with each replica have different state. Eventually consistent is requirement that once partitioning is resolved, last step, identified inconsistent copies of data items and updates to new version of data that is meaningful to the application. Application decides how to resolve the inconsistency not database.

## \* Distributed Query Processing

- Centralized system - criterion of cost is disk accesses while distributed system - Data transmission over the network and potential gain in performance from having several sites process parts of the query in parallel.

### \* Query Transformation:

- If no replicas are fragmented, we choose the replica which transmission cost is lowest. However, if replica is fragmented the choice is not so easy to make, since several joins or unions require to reconstruct the "account" relation

Ex =  $\sigma_{\text{branchname} = \text{"Hillside}} (\text{account})$ , Here  $\text{account} = \text{account1} \cup \text{account2}$

$$\sigma_{\text{branchname} = \text{"Hillside}} (\text{account1}) \cup \sigma_{\text{branchname} = \text{"Hillside}} (\text{account2})$$

### \* Simple Join Processing

- account  $\bowtie$  depositor  $\bowtie$  branch

$s_1$        $s_2$        $s_3$

strategies ① ship copies of all three relations to site  $s_1$  ↗ drawback

②  $s_1$  to  $s_2$  compute account  $\bowtie$  depositor  $\rightarrow$  temp

$s_2$  to  $s_3$  compute temp  $\bowtie$  branch

\* Join strategies that exploit Parallelism

$r_1 \bowtie r_2 \bowtie r_3 \bowtie r_4$  - output at  $s_1$

① ship  $r_1$  from  $s_1$  to  $s_2$  ② Compute  $r_{21} \bowtie r_{22}$  at  $s_2$  = temp<sub>1</sub>

③ simultaneously ship  $r_3$  from  $s_3$  to  $s_4$  ④ Compute  $r_{31} \bowtie r_{32}$  at  $s_4$  = temp<sub>2</sub>

⑤ ship temp<sub>1</sub> from  $s_2$  to  $s_1$  ⑥ ship temp<sub>2</sub> from  $s_4$  to  $s_1$

⑦ temp<sub>1</sub>  $\bowtie$  temp<sub>2</sub> at  $s_1$

more computation

more network traffic

(14) 5

## \* Semijoin strategy ( $r_1 \bowtie r_2$ )

→ Suppose expression  $r_1 \bowtie r_2$  where  $r_1$  is stored at S1 and  $r_2$  is stored at S2.

→ Problem: all tuples from  $r_2$  moved to S1 and then perform join, Network cost is too much

Let the schemas of  $r_1$  &  $r_2$  is R1 & R2

बिपीनकुमार अम. पटेल  
राष्ट्रीय जीवन बीमा निगम  
LIFE INSURANCE CORPORATION OF INDIA  
संस्था - अधिकारी के लिये सेविंग प्रोडक्शन



Bipinkumar M. Patel  
Member of the Zonal Manager's Club for Agents  
828 83A

① Compute  $\text{temp}_1 \leftarrow \pi_{R_1 \cap R_2}(r_1)$  at S1

② Ship  $\text{temp}_1$  from S1 to S2

③ Compute  $\text{temp}_2 \leftarrow r_2 \Delta \text{temp}_1$  at S2

④ Ship  $\text{temp}_2$  from S2 to S1

⑤ Compute  $r_1 \bowtie r_2 \Delta \text{temp}_2$  at S1

$$r_1 \bowtie r_2 \Delta \pi_{R_1 \cap R_2}(r_1) = (r_1 \Delta \pi_{R_1 \cap R_2}(r_1)) \bowtie r_2 = r_1 \bowtie r_2$$

→ Overhead of shipping  $\text{temp}_1$  will be dominated by the saving of shipping only a fraction of the tuples in  $r_2$ .

→  $\text{temp}_2$  should have less tuples than and then it strategy is useful

## \* Cloud Based Databases

### ① Data Representation

- data storage support XML & JSON for representing such complex data.
- Both data ~~structure~~ structure is flexible in set of attributes that a record contains
- store data with associated key and retrieve data with that key.
- In social networking application when someone post photo then it will send to his friends so whenever their friends open social media they received posts. "no delay"
- Cloud data storage systems are based on two primitive functions `put(key, value)` and `get(key)`
- In Bigtable, a record is not stored ~~as~~ single value but is instead split into component attributes that stored separately.
- Each attribute value is just a string as far as Bigtable concerned.
- Record identifier is used to fetch the record.
- Single instance of Bigtable can store data for multiple applications with multiple table per application, by simply prefixing the application name and table name to the record identifier.
- Data storage system typically allow multiple version of data items to be stored. Versions are often identified by timestamp, or integer value.

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता.जी.आर०-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.

- Lookups can simply pick data which has highest version.

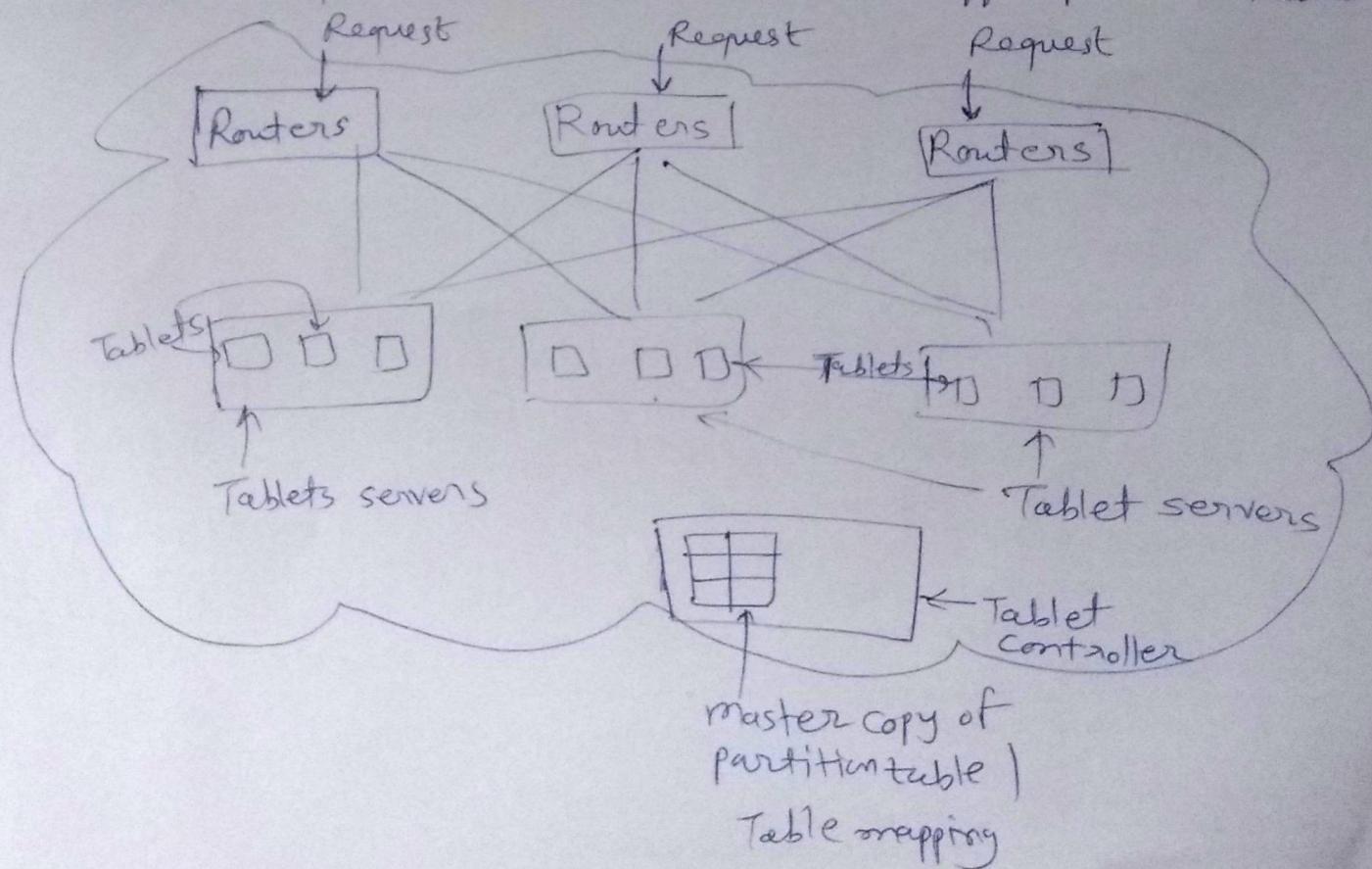
Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail: bipinpatellic@yahoo.com

## Partitioning and Retrieving Data

- Problems extremely large scale in data storage system, load ↑ on one site
- = Partition data relatively small units called tablets
- each tablet is a fragment of a table
  - The partitioning of data should be done on the search key, so that a request for specific key value is directed to a single tablet; otherwise each request would require processing at multiple sites, increasing the load on the system.
  - Range or hash partitioning is used
  - A site to which a tablet is assigned acts as the master site for that tablet. All updates are routed through this site, and updates are then propagated to replicas of the tablet.
  - Partitioning of data into tablets is dynamic. As data are inserted, if a tablet grows too big, it is broken into smaller parts.
  - even tablet is not large enough to merit being broken into smaller tablets and distributes over sites.
  - Controller site - Tracks the partition function to map a "get()" request to one or more tablets, and a mapping function from tablets to sites, to find which site were responsible for which tablet.
  - single tablet controller site would soon get overloaded. Instead, the mapping information replicated on a set of router sites, which route requests to the site with the appropriate tablet.



19  
⑥

## \* Transactions and Replication

- Such system support transactions on data within a single tablet, which controlled by a single master site.

- It update data item to be conditional on the current version number of the data item is more recent than the specified version number.

बिपीनकुमार अम. पटेल

LIC  
भारतीय जीवन बीमा निगम  
LIFE INSURANCE CORPORATION OF INDIA



संस्था - अधिकारीजी के लिये होर्डिंग प्रबोक क्लब

Bipinkumar M. Patel  
Member of the Zonal Manager's Club for Agents  
82883A

In system thousands of sites, at any time it is almost guaranteed that several of the sites will be down. A data storage system on the cloud must be able to continue normal processing even with many sites down since each table is controlled by a single master site, if the site fails the tablet should be reassigned to a different site that has a copy of the tablet, which becomes the new master site for tablet. When site fail new controller use log and updates, lookups can be performed on the tablet.

Replication at a remote site is for high availability propagating updates to secondary indices. Secondary indices are basically tables, partitioned just like regular tables, based on the index search key, an update of record in a table can be mapped to updates of one or more tablets in a secondary index on the table.

## \* Traditional Databases on Cloud

Virtual machines - user was given the illusion of having private computer. It provides great flexibility. User can use their own O.S. & softwares. Entire computer can be allocated to each virtual machine of a client whose virtual machine have a high load. Easy to add or subtract computing power as workloads grow shrink. Applications are easily parallelized.

## \* Challenges with Cloud-Based Databases

Distributed database system require frequent communication and coordination among sites for

① access to data on another machine

② obtaining lock on remote data

③ ensuring atomic transaction commit via 2-phase Commit.

Effective query optimization requires that the optimizer have accurate cost measures for operations.

Lacking knowledge of physical placement of data, the optimizer has to rely on estimates that may be inaccurate.

कार्यालय : परबड़ी के पास, पो. समरखा.  
निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणंद-388360.  
फोन : (नि) (02692) 256673 (मो) 98242 32755.

- Replication is done by vendor without knowledge of application e.g. It is controlled by cloud not by database.

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. / P. (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

- Now a days, application itself need to take responsibility for ~~consistency~~  
<sup>only</sup>
- Users of cloud computing have security and legal liability problem because data held by another organization.
- Main problem, when cloud vendor store data in foreign country  
e.g. if a german company store data in "NY" then U.S. government can see that data without client permission.

## \* Directory Systems

(person)

- Directory is a listing of information about some class of objects
- It can be used to find information about specific object
- Directories
  - ↳ lookups in forward direction - white pages
  - ↳ lookups in reverse direction - yellow pages

## \* Directory Access Protocols

- Several protocols are developed to provide a standardized way of accessing data in directory. Most used, is "Lightweight Directory Access Protocol"
- why JDBC/ODBC not used?

- ① This protocol is simplified, for a limited type of data access, evolved parallel to ODBC/JDBC.
- ② Directory system provides nice hierarchical naming mechanism
  - Data can be partitioned amongst multiple server for different parts of hierarchy.
  - directory system can be setup to automatically forward queries made at one site to other site, without user intervention.

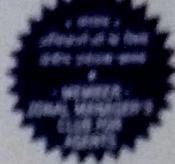
## \* LDAP Data model

- Directory store entries. Each entry must have distinguished name (DN)
- A DN is in turn made up of sequence of relative distinguished names (RDNs)
- DN is combination of attributes
- The set of RDNs for a DN is defined by schema of the directory system.
- It provides string, binary, time types, tel ← for telephone
- all attributes are multivalued by default.
- LDAP provides object classes with names and types attribute
- Entries are organized into a directory information tree, according to their distinguished names. Entries at leaf level represent specific objects. Internal nodes represent - organization Unit, country. Children of a node inherit the DN of the parent, and add RDNs e.g. internal node with DN C=USA child node have DN starting with C=USA & further RDNs such as "O" or "OU"
- DN of an entry can be generated by traversing path from root as "O" or "OU".
- leaf level can be pointing another entry
- Entries can have more than one DN

[read from PPT]

## CH-22 Complex Data types

①



- Recently, complex data types such as addresses work fine because address could be viewed as an atomic data item and it is easy to store data.

### Complex type system

- Composite attributes, multivalued
- generalization, specialization

बिपीनकुमार अम. पटेल

पद्धति - अधिकारी के लिए दोस्री प्रकाशन वर्ष

Bipinkumar M. Patel  
Member of the Zonal Manager's Club for Agents  
828 83A

Rather than view a database as a set of records, users of applications view it as a set of objects, without complex data types translation to relational model.

- These objects require set of records for representation
- Easy-to-use interface requires a one-to-one correspondence between the user's intuitive notion of an object and the database system's notion of a data item.

Ex Library application

Authors - book may have list of authors, which can represent as an array

Keywords - keywords include one or more specified keywords

Publisher - Publisher consisting name and branch.

Above all three fields are monatomic which mean they are no-INF.

Same system can implement in different way and have 4NF schema whereas non-INF design makes many types of queries easier.

\* Structured Types and Inheritance in SQL

\* Structured Types: user define type

Create type Name as (firstname varchar(20), lastname varchar(20))

Create type Address as (street varchar(20), city varchar(20)) final;

We can now create composite attribute in relation

Create table person(name Name, address Address, Dob date);

Components of composite attribute can accessed using ".":

E.g. name.firstname

We can also create a table whose rows are of user define type

Ex Create type personType as (name Name, address Address, Dob date);

Create table person of personType;

\* [In other way]

Create table person (name raw (firstname  
, address raw (street varchar(20),  
city varchar(20)),  
Dob date));

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. बी.आर०-388360.

फोन : (रि) (02692) 256673 (मो) 98242 32755.

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

Here name, and address are unnamed type and row of table also mon...

structured type can have methods define on it.

```

create type PersonType as (name Name, address Address, DOB date)
    not final
    method age (Date date)
        returns interval year;
create instance method age (Date date) returns interval year;
begin
    for PersonType
    return Date - self.DOB;
end

```

→ This method executes on an instance of the personType.  
 → The variable "self" refers to the person instance on which the method is invoked.

\* After **Person** table created, fetch age and last-name

**Select** name, lastname, age (current\_date) from person;

\* **create** constructor for structured Type

```

create function Name (firstname varchar(20), lastname varchar(20))
returns Name
begin
    set self.firstname = firstname
    set self.lastname = lastname
end

```

→ By default every structured type has a constructor with no arguments  
 - There can be more than one constructor for same structured type although they distinguishable by the number of arguments and types of their arguments

insert into person values (new Name ('mit', 'Patel'), Address ("main", "anand"), date '1999-08-12');

## \* Type Inheritance



- Student and teacher also persons so both have same attribute hence person type can be inherited

Create type Student under Person (degree nvarchar(20))

Create type Teacher under Person (salary integer);

बिपीनकुमार अम. पटेल

मरवाड़ - अभियांत्रियों के लिये शैक्षण्य प्रबंधक वर्षा

Bipinkumar M. Patel  
Member of the Zonal Manager's Club for Agents  
82883A

- Here methods of structured type are inherited by its subtypes.

- declaring method can be override using "overriding method"
- final says that subtype may not be created nonfinal subtype may be created

\* Here we can inherit student and teacher both simultaneously

\* If student & Teacher both have department attribute than both but they are not cause any problem because of both are inherited from Create type Teching Assistant under Student, Teacher

each value must be associated with one specific type called

## \* Table inheritance

Create table people of Person;

Create table student of Student under people;

Create table teacher of Teacher under people;

→ all person's attributes inherited into teacher and student table.

Delete from people where name = "mit"

→ This will affect all three table because people is supertable and student and teacher are subtable - all tuples will be deleted

→ Delete from only people where name = "mit"

- only people will alter

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणंद-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.

→ multiple inheritance of table is not supported by SQL

e.g. Create table x of x1  
under P1, P2

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

- tuples in a subtable and parent table correspond if they have the same values for all inherited attributes. Thus, corresponding tuples represent the same entity.
- \* The consistency require for subtables are
  - ① Each tuple of supertable can correspond to at most one tuple in each of its immediate subtables.
  - ② all the tuples corresponding to each other must be derived from one tuple.

→ without first condition two tuples in students

- since SQL does not support multiple inheritance, it prevents a person from being both a teacher and a student.

### \* Array and Multiset Types in SQL

- Create type Publisher as (name varchar(20), branch varchar(20));  
Create type Book as (title varchar(20), author varchar(20) array[10],  
Create table b of Book;  
publish Publisher, Keyword varchar(20) multiset);

### \* Creating and accessing Collection Values

insert into b values ('Compiler', array['mit', 'smith'], new Publisher('me', 'NY'),  
 → Assume Publisher constructor multiset['parsing', 'analysis']);

### \* Parsing Collection-valued Attributes

Select title from b where 'Parsing' in unnest(Keyword));

Select title, A.author, A.position from b, unnest(b.author)  
 → with ordinality as A (author, position);

\* Nesting and Unnesting extra attribute return - position of the element in the array

Select title, A.author, publisher.name as pub-name, publisher.branch as pub-branch,  
 ↑ K.keyword from b, unnest(b.author) as A (author),  
 output is INF ↓ nest unnest (b.Keyword) as K (Keyword);

Select title, author, Publisher (pub-name, pub-branch), collect (Keyword)  
 from books group by title, author, publisher

③ Select title, array (select author as A where A.title = B.title), Publisher (pub-name, pub-name),  
 multiset (select Keyword from Keyword as K where K.title = B.title)  
 from books as B

## Ch-22 \* Object-Identity and Reference Types in SQL



भारतीय जीवन स्थिरा निगम  
LIFE INSURANCE CORPORATION OF INDIA

Create type Department (name Varchar(20),  
 head ref(Person) scope people);

- Here scope works like foreign keys. We can omit scope people  
Create table department of Department (head with options  
 scope people);

\* self-referential attribute

Create table people of person  
 ref is personid system generated;



बिपीनकुमार अम. पटेल  
सदस्य - अधिकारी के लिए हैलिंग प्राइवेट कल्पना

Bipinkumar M. Patel  
Member of the Zonal Manager's Club for Agents  
828 83A

In order to initialize a reference attribute, we need to get the identifier of the tuple that is to be referred. We may first create the tuple with a null reference and then set the reference separately:

insert into department value ('cs', null);

update department set head = (select p.personid from people as p  
 where name = 'cs' ; where name = 'mit')

Alternative to system generated identifiers is to allow users to generate  
Create type Person (name Varchar(20), address Varchar(20)) ref using Varchar(20);  
Create table people of Person ref is personid user generated;

insert into people values (personid, name, address)  
 values ('123', 'mit', 'NY');

insert into department values ('cs', '123');

→ It is possible to use an existing primary key value as the identifier, by including the ref from clause in the type definition:

Create type Person (name Varchar(20) primary key, address Varchar(20))  
 ref from (name);

Create table people of Person ref is personid derived;  
 insert into department values ('cs', 'mit');

→ Reference are dereferenced by the "→" symbol

Select head → name, head → address from department;

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणं-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.

→ "dref" returning the tuple pointed by a reference and then access its attributes.

Select dref(head).name  
 from department;

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

## Object-Relational Mapping

- It is built on top of a traditional relational database, and allows a programmer to define a mapping between tuples in database relations and objects in the programming language.
- An object, or set of objects, can be retrieved based on a selection condition on its attributes; relevant data are retrieved from the underlying database based on the selection conditions.
- Mapping from objects to relations is used to correspondingly update, insert or delete tuples in the database.
- Main goal of mapping is proving object model, benefits of using a robust relational database underneath.
- This system provide significant performance gains over direct access to the underlying database.
- This mapping also provide query languages that allow to ~~not~~ write query for direct object Model; such queries are translated into SQL and result object created from SQL query results.
- It suffer from significant overheads for bulk database updates, and may provide only limited querying capability.
- Benefits of mapping exceed the drawbacks for many applications.

## \* Hadoop



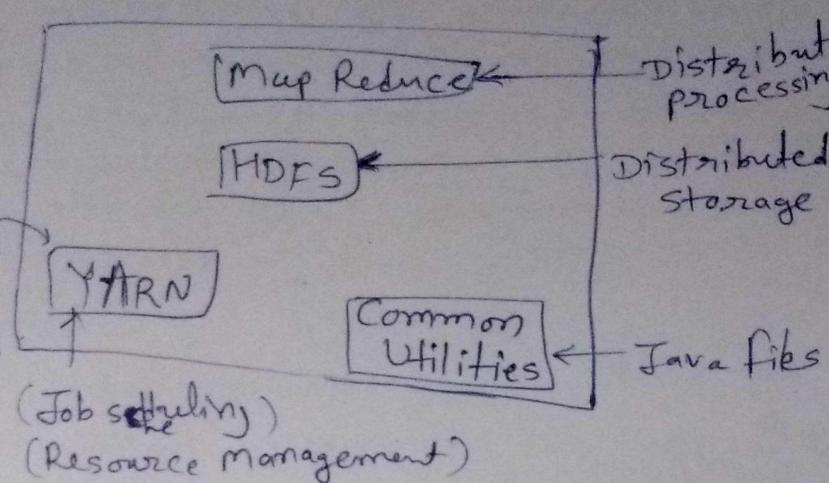
भारतीय जीवन बीमा नियम  
LIFE INSURANCE CORPORATION OF INDIA

Yet another Resource Negotiator  
बिपीनकुमार अम. पटेल

सरदार - अधिकारी के लिए होटिंग प्रशंसक कल्याण

Bipinkumar M. Patel

Member of the Zonal Manager's Club for Agents  
828 83A



- Hadoop follows a master-slave architecture using HDFS and Map-Reduce respectively

- A file on HDFS is split into multiple blocks and each is replicated within the hadoop cluster. A default block size is 64 MB. It can be extended up to 256 MB

- Namenode and DataNode are two critical components of the Hadoop HDFS architecture.

- Application data is stored on DataNodes & File system metadata is stored on ~~NameNode~~ NameNode.

- The NameNode and DataNode communicate with each other using TCP based protocols.

- HDFS must satisfy prerequisites-

- ① All hard drives should have a high throughput
- ② Good network speed to manage intermediate data transfer and block replications.

### \* NameNode

- all the files and directories in namespace are represented by Inodes that contain various attributes like permissions, modification timestamp, disk space quota and access times.

- Two files fsimage and edits are used for persistence

[fsimage] file contains the Inodes and list of block which define the metadata. It has complete snapshot of the file systems.

[edit] file contains any modifications that have been performed on the content of the fsimage file.

- Create new fsimage snapshot everytime the namespace is being altered.

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता.जी.आणंद-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.

- When NameNode starts, fsimage file is loaded and then the contents of the edit file are applied to recover the latest state of the file system.

Problem is file grows and consume more space, Slow down running.

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

- If the hadoop cluster has not been restarted for months together then there will be a huge downtime as the size of the edits file will be increase. This is when Secondary NameNode comes to the rescue. Secondary NameNode gets the fsimage and edits log from the primary NameNode at regular intervals and loads both the fsimage and edit logs file to the main memory by applying each operation from edits log file to fsimage.

### \* DataNode

- A DataNode can perform CPU intensive jobs like semantic and language analysis, statistics and machine learning tasks and indexing. A DataNode needs lot of I/O for data processing and transfer.
- On startup every DataNode connects to the NameNode and performs a handshake to verify the Namespace ID and software version of the DataNode. If either of them does not match dataNode shutdown.
- DataNode sends heartbeat to the NameNode every 3 seconds to confirm that the dataNode is available.

### \* MapReduce

- The execution of a MapReduce job begins when the client submits the job configuration to the Job Tracker that specifies the map, combine and reduce functions along with the location for input and output data.
- Job tracker identifies the number of splits based on the input path and selected task trackers based on their network vicinity to the data sources. Job tracker send signal to selected Task trackers.
- After the completion of work task trackers notify job tracker.
- When all task trackers are done, job tracker collect all work

### Hadoop

→ Used for structured, semistructured and unstructured data	RDBMS
→ Use for large dataset	→ mainly for structured data
→ High query language	→ Average size dataset
→ Required schema on read (dynamic)	→ Structured query language
→ Read & writes are fast	→ required schema on write (static)
→ free	→ Reads are fast
→ Throughput High	→ License
→ Work on key/pair (Data objects) value	→ Throughput low
→ Horizontal scalability	→ Works on Relational Tables
→ Integrity low	→ Vertical scalability
→ Commodity/utility hardware	→ Integrity High
→ Analytics, Data Discovery	→ High end servers
	→ online transaction processing