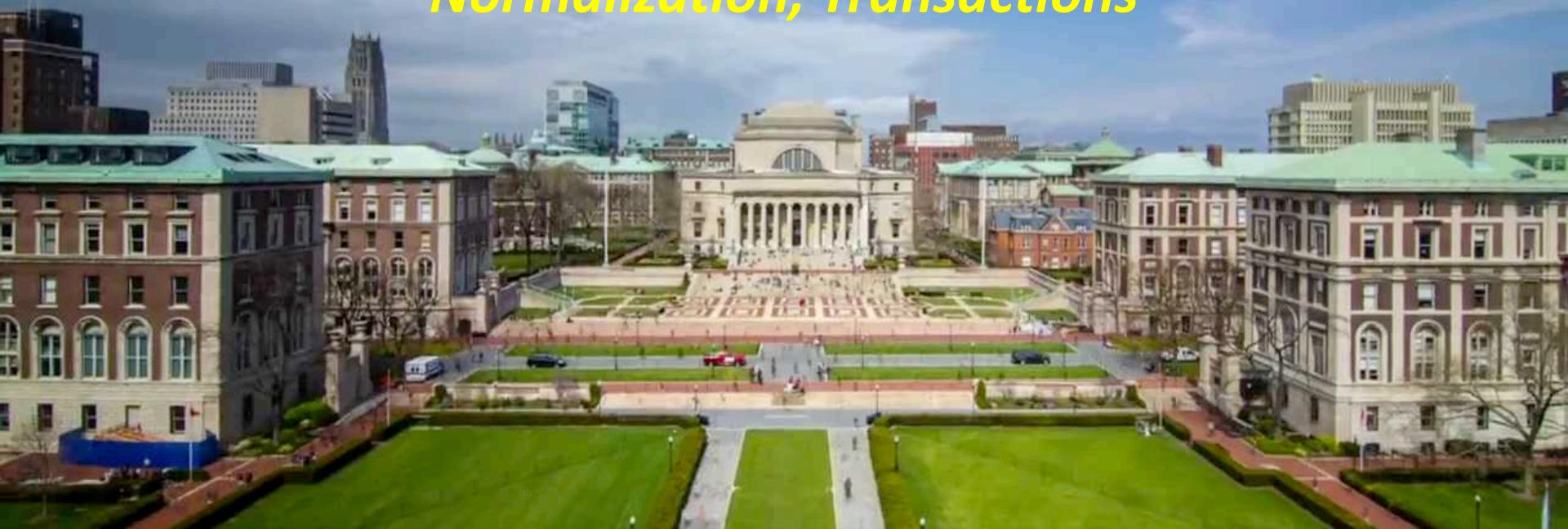


*Lecture 11:*  
*Normalization, Transactions*



*Lecture 11:  
Normalization, Transactions*

We will start in a couple of minutes.

# Contents

# Contents

# Normalization



# Chapter 7: Normalization

Database System Concepts, 7<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Outline

- Features of Good Relational Design
- Functional Dependencies
- Decomposition Using Functional Dependencies
- Normal Forms
- Functional Dependency Theory
- ~~Algorithms for Decomposition using Functional Dependencies~~
- ~~Decomposition Using Multivalued Dependencies~~
- ~~More Normal Form~~
- Atomic Domains and First Normal Form
- ~~Database Design Process~~
- ~~Modeling Temporal Data~~



# Overview of Normalization



# Features of Good Relational Designs

- Hypothetically, assume our schema originally had one relation that provided information about faculty and departments.
- Suppose we combine *instructor* and *department* into *in\_dep*, which represents the natural join on the relations *instructor* and *department*

<i>ID</i>	<i>name</i>	<i>salary</i>	<i>dept_name</i>	<i>building</i>	<i>budget</i>
22222	Einstein	95000	Physics	Watson	70000
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000
83821	Brandt	92000	Comp. Sci.	Taylor	100000
15151	Mozart	40000	Music	Packard	80000
33456	Gold	87000	Physics	Watson	70000
76543	Singh	80000	Finance	Painter	120000

- There is repetition of information
- Need to use null values (if we add a new department with no instructors)



# Decomposition

- The only way to avoid the repetition-of-information problem in the *in\_dep* schema is to decompose it into two schemas – instructor and *department* schemas.
- Not all decompositions are good. Suppose we decompose

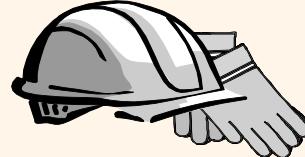
*employee*(*ID*, *name*, *street*, *city*, *salary*)

into

*employee1* (*ID*, *name*)

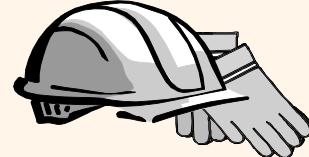
*employee2* (*name*, *street*, *city*, *salary*)

- The problem arises when we have two employees with the same name
- The next slide shows how we lose information -- we cannot reconstruct the original *employee* relation -- and so, this is a **lossy decomposition**.



# The Evils of Redundancy

- ❖ *Redundancy* is at the root of several problems associated with relational schemas:
  - redundant storage, insert/delete/update anomalies
- ❖ Integrity constraints, in particular *functional dependencies*, can be used to identify schemas with such problems and to suggest refinements.
- ❖ Main refinement technique: *decomposition* (replacing ABCD with, say, AB and BCD, or ACD and ABD).
- ❖ Decomposition should be used judiciously:
  - Is there reason to decompose a relation?
  - What problems (if any) does the decomposition cause?



## Example (Contd.)

- ❖ Problems due to  $R \rightarrow W$  :
  - Update anomaly: Can we change  $W$  in just the 1st tuple of SNLRWH?
  - Insertion anomaly: What if we want to insert an employee and don't know the hourly wage for his rating?
  - Deletion anomaly: If we delete all employees with rating 5, we lose the information about the wage for rating 5!

Will 2 smaller tables be better?

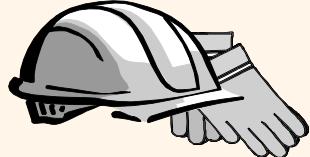
Wages

R	W
8	10
5	7

Hourly\_Emps2

S	N	L	R	H
123-22-3666	Attishoo	48	8	40
231-31-5368	Smiley	22	8	30
131-24-3650	Smethurst	35	5	30
434-26-3751	Guldu	35	5	32
612-67-4134	Madayan	35	8	40

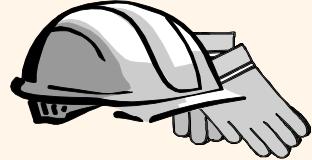
S	N	L	R	W	H
123-22-3666	Attishoo	48	8	10	40
231-31-5368	Smiley	22	8	10	30
131-24-3650	Smethurst	35	5	7	30
434-26-3751	Guldu	35	5	7	32
612-67-4134	Madayan	35	8	10	40



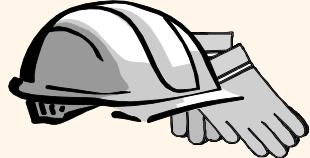
# Functional Dependencies (FDs)

- ❖ A functional dependency  $X \rightarrow Y$  holds over relation  $R$  if, for every allowable instance  $r$  of  $R$ :
  - $t1 \in r, t2 \in r, \pi_X(t1) = \pi_X(t2)$  implies  $\pi_Y(t1) = \pi_Y(t2)$
  - i.e., given two tuples in  $r$ , if the  $X$  values agree, then the  $Y$  values must also agree. ( $X$  and  $Y$  are *sets* of attributes.)
- ❖ An FD is a statement about *all* allowable relations.
  - Must be identified based on semantics of application.
  - Given some allowable instance  $r1$  of  $R$ , we can check if it violates some FD  $f$ , but we cannot tell if  $f$  holds over  $R$ !
- ❖  $K$  is a candidate key for  $R$  means that  $K \rightarrow R$ 
  - However,  $K \rightarrow R$  does not require  $K$  to be *minimal*!

## *Example: Constraints on Entity Set*



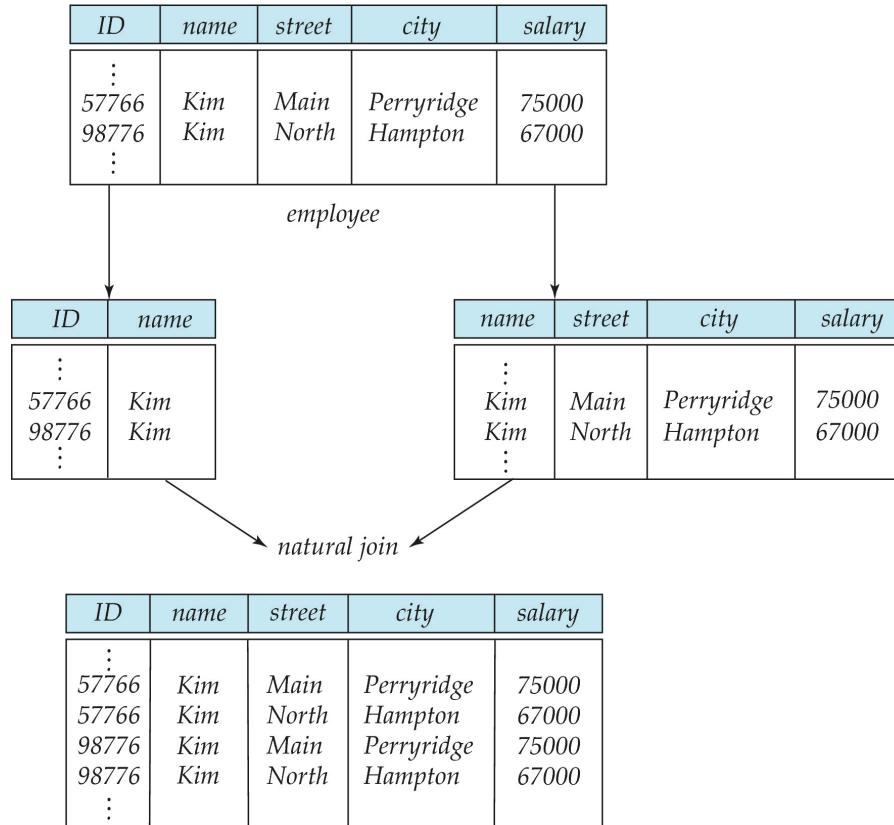
- ❖ Consider relation obtained from Hourly\_Emps:
  - Hourly\_Emps (ssn, name, lot, rating, hrly\_wages, hrs\_worked)
- ❖ Notation: We will denote this relation schema by listing the attributes: **SNLRWH**
  - This is really the *set* of attributes {S,N,L,R,W,H}.
  - Sometimes, we will refer to all attributes of a relation by using the relation name. (e.g., Hourly\_Emps for SNLRWH)
- ❖ Some FDs on Hourly\_Emps:
  - *ssn is the key:*  $S \rightarrow \text{SNLRWH}$
  - *rating determines hrly\_wages:*  $R \rightarrow W$

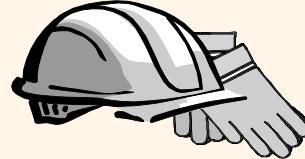


# Reasoning About FDs

- ❖ Given some FDs, we can usually infer additional FDs:
  - $ssn \rightarrow did$ ,  $did \rightarrow lot$  implies  $ssn \rightarrow lot$
- ❖ An FD  $f$  is *implied by* a set of FDs  $F$  if  $f$  holds whenever all FDs in  $F$  hold.
  - $F^+ = \text{closure of } F$  is the set of all FDs that are implied by  $F$ .
- ❖ Armstrong's Axioms (X, Y, Z are sets of attributes):
  - Reflexivity: If  $X \subseteq Y$ , then  $Y \rightarrow X$
  - Augmentation: If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$  for any  $Z$
  - Transitivity: If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$
- ❖ These are *sound* and *complete* inference rules for FDs!

# Lossy Decomposition





## Lossless Decomposition

- ❖ Let  $R$  be a relation schema and let  $R_1$  and  $R_2$  form a decomposition of  $R$ .
  - That is  $R = R_1 \cup R_2$
  - *Note: This notation is confusing. This is a statement about the schema, not about the data in relations.*
- ❖ We say that the decomposition is a **lossless decomposition** if there is no loss of information by replacing  $R$  with the two relation schemas  $R_1 \cup R_2$
- ❖ Formally,

$$\Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$$

- ❖ And, conversely a decomposition is lossy if

$$r \subset \Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$$



# Normalization Theory

- Decide whether a particular relation  $R$  is in “good” form.
- In the case that a relation  $R$  is not in “good” form, decompose it into set of relations  $\{R_1, R_2, \dots, R_n\}$  such that
  - Each relation is in good form
  - The decomposition is a lossless decomposition
- Our theory is based on:
  - Functional dependencies
  - Multivalued dependencies



# Functional Dependencies

- There are usually a variety of constraints (rules) on the data in the real world.
- For example, some of the constraints that are expected to hold in a university database are:
  - Students and instructors are uniquely identified by their ID.
  - Each student and instructor has only one name.
  - Each instructor and student is (primarily) associated with only one department.
  - Each department has only one value for its budget, and only one associated building.



# Functional Dependencies (Cont.)

- An instance of a relation that satisfies all such real-world constraints is called a **legal instance** of the relation;
- A legal instance of a database is one where all the relation instances are legal instances
- Constraints on the set of legal relations.
- Require that the value for a certain set of attributes determines uniquely the value for another set of attributes.
- A functional dependency is a generalization of the notion of a *key*.



# Functional Dependencies Definition

- Let  $R$  be a relation schema

$$\alpha \subseteq R \text{ and } \beta \subseteq R$$

- The **functional dependency**

$$\alpha \rightarrow \beta$$

**holds on  $R$**  if and only if for any legal relations  $r(R)$ , whenever any two tuples  $t_1$  and  $t_2$  of  $r$  agree on the attributes  $\alpha$ , they also agree on the attributes  $\beta$ . That is,

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

- Example: Consider  $r(A,B)$  with the following instance of  $r$ .

1	4
1	5
3	7

- On this instance,  $B \rightarrow A$  hold;  $A \rightarrow B$  does **NOT** hold,



# Closure of a Set of Functional Dependencies

- Given a set  $F$  set of functional dependencies, there are certain other functional dependencies that are logically implied by  $F$ .
  - If  $A \rightarrow B$  and  $B \rightarrow C$ , then we can infer that  $A \rightarrow C$
  - etc.
- The set of **all** functional dependencies logically implied by  $F$  is the **closure** of  $F$ .
- We denote the *closure* of  $F$  by  $F^+$ .



# Keys and Functional Dependencies

- $K$  is a superkey for relation schema  $R$  if and only if  $K \rightarrow R$
- $K$  is a candidate key for  $R$  if and only if
  - $K \rightarrow R$ , and
  - for no  $\alpha \subset K$ ,  $\alpha \rightarrow R$
- Functional dependencies allow us to express constraints that cannot be expressed using superkeys. Consider the schema:

*in\_dep (ID, name, salary, dept\_name, building, budget ).*

We expect these functional dependencies to hold:

$dept\_name \rightarrow building$

$ID \rightarrow building$

but would not expect the following to hold:

$dept\_name \rightarrow salary$

DFF Note:

- In the current data:  $ID \rightarrow dept\_name \rightarrow building$
- But
- In the schema,  $dept\_name$  may be NULL..



# Use of Functional Dependencies

- We use functional dependencies to:
  - To test relations to see if they are legal under a given set of functional dependencies.
    - ▶ If a relation  $r$  is legal under a set  $F$  of functional dependencies, we say that  $r$  **satisfies**  $F$ .
  - To specify constraints on the set of legal relations
    - ▶ We say that  $F$  **holds on**  $R$  if all legal relations on  $R$  satisfy the set of functional dependencies  $F$ .
- Note: A specific instance of a relation schema may satisfy a functional dependency even if the functional dependency does not hold on all legal instances.
  - For example, a specific instance of *instructor* may, by chance, satisfy  $name \rightarrow ID$ .



# Normal Forms



# Boyce-Codd Normal Form

- A relation schema  $R$  is in BCNF with respect to a set  $F$  of functional dependencies if for all functional dependencies in  $F^+$  of the form

$$\alpha \rightarrow \beta$$

where  $\alpha \subseteq R$  and  $\beta \subseteq R$ , at least one of the following holds:

- $\alpha \rightarrow \beta$  is trivial (i.e.,  $\beta \subseteq \alpha$ )
- $\alpha$  is a superkey for  $R$

## DFF Note:

- The theoretical treatment is no conveying the practical intent.
- If  $\alpha$  is a superkey, I can set a primary key/unique constraint on  $\alpha$ .
- Consider tables with address info in the rows.



# Decomposing a Schema into BCNF

- Let  $R$  be a schema  $R$  that is not in BCNF. Let  $\alpha \rightarrow \beta$  be the FD that causes a violation of BCNF.
- We decompose  $R$  into:
  - $(\alpha \cup \beta)$
  - $(R - (\beta - \alpha))$
- In our example of *in\_dep*,
  - $\alpha = \text{dept\_name}$
  - $\beta = \text{building}, \text{budget}$and *in\_dep* is replaced by
  - $(\alpha \cup \beta) = (\text{dept\_name}, \text{building}, \text{budget})$
  - $(R - (\beta - \alpha)) = (\text{ID}, \text{name}, \text{dept\_name}, \text{salary})$

DFF Note – again this is baffling

- $R = (\text{id}, \text{name\_last}, \text{name\_first}, \text{street}, \text{city}, \text{state}, \text{zipcode})$
- $\alpha \rightarrow \beta$  means  $\text{zipcode} \rightarrow (\text{city}, \text{state})$
- $(\alpha \cup \beta) = (\text{zipcode}, \text{city}, \text{state})$ , which we call Address
- $R - (\beta - \alpha) = R - \beta + \alpha = (\text{id}, \text{name\_last}, \text{name\_first}, \text{zipcode})$ , which we call Person
- Setting primary key ID in Address and zipcode on Address preserves the dependency and is lossless.

Note:

- This is an example only.
- Zip code does not imply city or state in real world.



# BCNF and Dependency Preservation

- It is not always possible to achieve both BCNF and dependency preservation
- Consider a schema:  
$$\text{dept\_advisor}(s\_ID, i\_ID, \text{department\_name})$$
- With function dependencies:  
$$i\_ID \rightarrow \text{dept\_name}$$
  
$$s\_ID, \text{dept\_name} \rightarrow i\_ID$$
- $\text{dept\_advisor}$  is not in BCNF
  - $i\_ID$  is not a superkey.
- Any decomposition of  $\text{dept\_advisor}$  will not include all the attributes in  
$$s\_ID, \text{dept\_name} \rightarrow i\_ID$$
- Thus, the composition is NOT be dependency preserving



## Third Normal Form

- A relation schema  $R$  is in **third normal form (3NF)** if for all:

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- $\alpha \rightarrow \beta$  is trivial (i.e.,  $\beta \in \alpha$ )
- $\alpha$  is a superkey for  $R$
- Each attribute  $A$  in  $\beta - \alpha$  is contained in a candidate key for  $R$ .

(**NOTE**: each attribute may be in a different candidate key)

- If a relation is in BCNF it is in 3NF (since in BCNF one of the first two conditions above must hold).
- Third condition is a minimal relaxation of BCNF to ensure dependency preservation (will see why later).



## 3NF Example

- Consider a schema:

$dept\_advisor(s\_ID, i\_ID, dept\_name)$

- With function dependencies:

$i\_ID \rightarrow dept\_name$

$s\_ID, dept\_name \rightarrow i\_ID$

- Two candidate keys =  $\{s\_ID, dept\_name\}$ ,  $\{s\_ID, i\_ID\}$
- We have seen before that  $dept\_advisor$  is not in BCNF
- $R$ , however, is in 3NF
  - $s\_ID, dept\_name$  is a superkey
  - $i\_ID \rightarrow dept\_name$  and  $i\_ID$  is NOT a superkey, but:
    - ▶  $\{ dept\_name \} - \{ i\_ID \} = \{ dept\_name \}$  and
    - ▶  $dept\_name$  is contained in a candidate key



# Comparison of BCNF and 3NF

- Advantages to 3NF over BCNF. It is always possible to obtain a 3NF design without sacrificing losslessness or dependency preservation.
- Disadvantages to 3NF.
  - We may have to use null values to represent some of the possible meaningful relationships among data items.
  - There is the problem of repetition of information.

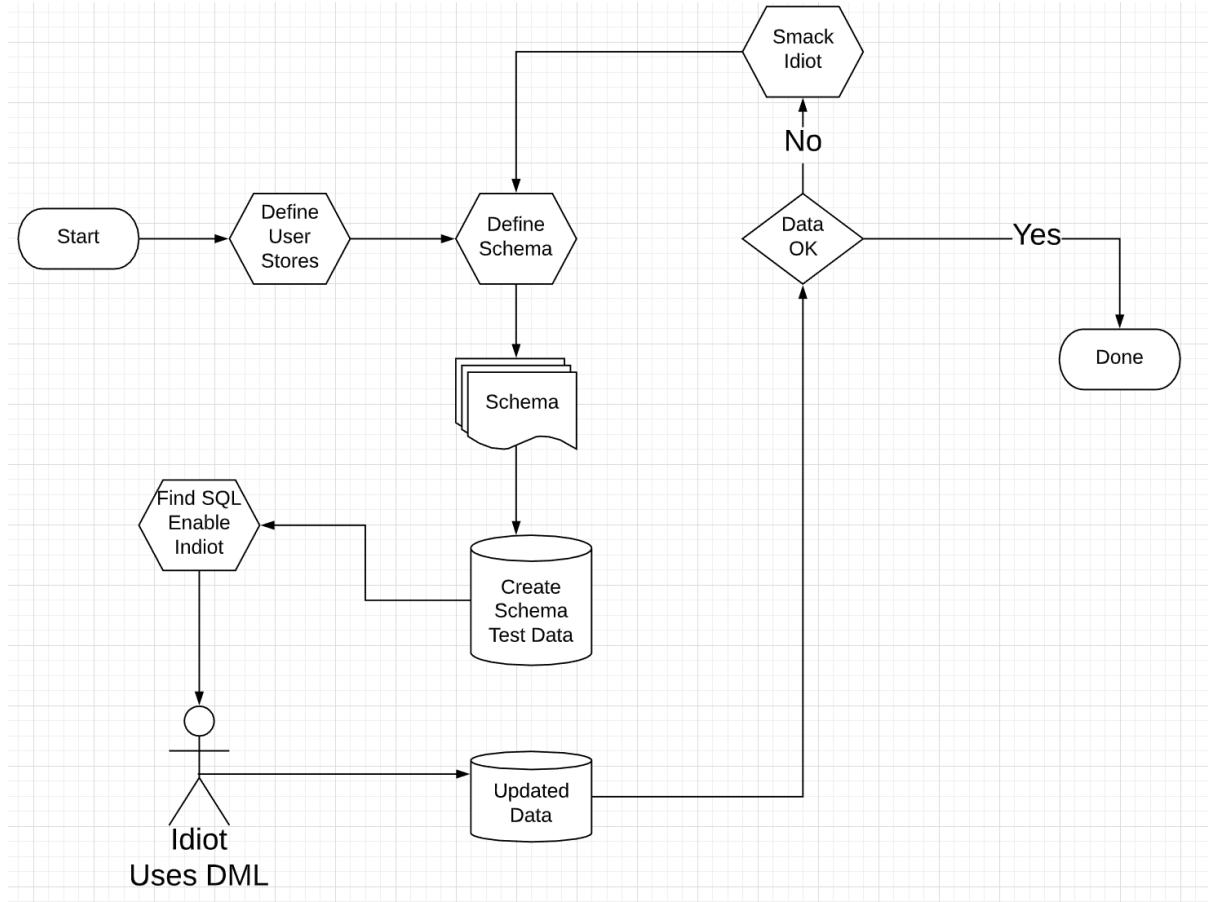


# Goals of Normalization

- Let  $R$  be a relation scheme with a set  $F$  of functional dependencies.
- Decide whether a relation scheme  $R$  is in “good” form.
- In the case that a relation scheme  $R$  is not in “good” form, need to decompose it into a set of relation scheme  $\{R_1, R_2, \dots, R_n\}$  such that:
  - Each relation scheme is in good form
  - The decomposition is a lossless decomposition
  - Preferably, the decomposition should be dependency preserving.

# **Universal – Ferguson's Law of Useful Normalization and Konstraints (U–FLUNK)**

# U-FLUNK is a Workflow Process





# Denormalization for Performance

- May want to use non-normalized schema for performance
- For example, displaying *prereqs* along with *course\_id*, and *title* requires join of *course* with *prereq*
- Alternative 1: Use denormalized relation containing attributes of *course* as well as *prereq* with all above attributes
  - faster lookup
  - extra space and extra execution time for updates
  - extra coding work for programmer and possibility of error in extra code
- Alternative 2: use a materialized view defined a  $course \bowtie prereq$ 
  - Benefits and drawbacks same as above, except no extra coding work for programmer and avoids possible errors



# First Normal Form

- Domain is **atomic** if its elements are considered to be indivisible units
  - Examples of non-atomic domains:
    - ▶ Set of names, composite attributes
    - ▶ Identification numbers like CS101 that can be broken up into parts
- A relational schema R is in **first normal form** if the domains of all attributes of R are atomic
- Non-atomic values complicate storage and encourage redundant (repeated) storage of data
  - Example: Set of accounts stored with each customer, and set of owners stored with each account
  - We assume all relations are in first normal form (and revisit this in Chapter 22: Object Based Databases)



# First Normal Form (Cont.)

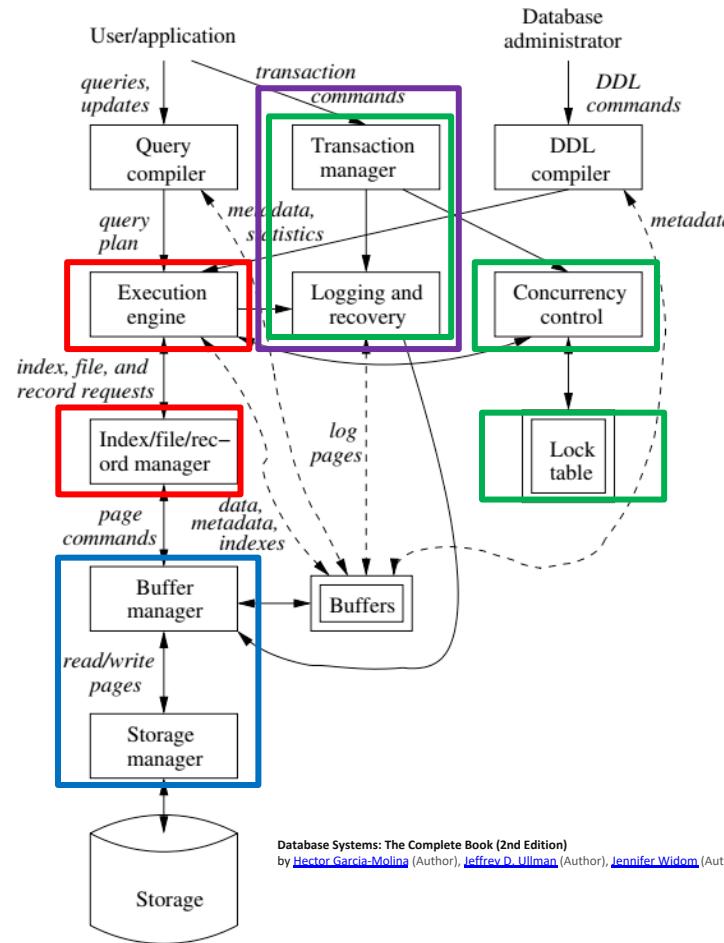
- Atomicity is actually a property of how the elements of the domain are used.
  - Example: Strings would normally be considered indivisible
  - Suppose that students are given roll numbers which are strings of the form *CS0012* or *EE1127*
  - If the first two characters are extracted to find the department, the domain of roll numbers is not atomic.
  - Doing so is a bad idea: leads to encoding of information in application program rather than in the database.

# Transactions Isolation Recovery

# Core Concepts

# Data Management

- Find things quickly.
- Load/Save quickly.
- Control access.
- Durability



Database Systems: The Complete Book (2nd Edition)  
by Hector Garcia-Molina (Author), Jeffrey D. Ullman (Author), Jennifer Widom (Author)

# Core Transaction Concept is ACID Properties

<http://slideplayer.com/slide/9307681>

## Atomic

“ALL OR NOTHING”

Transaction cannot be subdivided

## Consistent

Transaction → transforms database from one consistent state to another consistent state

**ACID**

## Isolated

Transactions execute independently of one another

Database changes not revealed to users until after transaction has completed

## Durable

Database changes are permanent  
The permanence of the database's consistent state

A *transaction* is a very small unit of a program and it may contain several low-level tasks. A transaction in a database system must maintain **Atomicity**, **Consistency**, **Isolation**, and **Durability** – commonly known as ACID properties – in order to ensure accuracy, completeness, and data integrity.

- **Atomicity** – This property states that a transaction must be treated as an atomic unit, that is, either all of its operations are executed or none. There must be no state in a database where a transaction is left partially completed. States should be defined either before the execution of the transaction or after the execution/abortion/failure of the transaction.
- **Consistency** – The database must remain in a consistent state after any transaction. No transaction should have any adverse effect on the data residing in the database. If the database was in a consistent state before the execution of a transaction, it must remain consistent after the execution of the transaction as well.
- **Durability** – The database should be durable enough to hold all its latest updates even if the system fails or restarts. If a transaction updates a chunk of data in a database and commits, then the database will hold the modified data. If a transaction commits but the system fails before the data could be written on to the disk, then that data will be updated once the system springs back into action.
- **Isolation** – In a database system where more than one transaction are being executed simultaneously and in parallel, the property of isolation states that all the transactions will be carried out and executed as if it is the only transaction in the system. No transaction will affect the existence of any other transaction.



# Transaction Concept

- A **transaction** is a *unit* of program execution that accesses and possibly updates various data items.
- E.g., transaction to transfer \$50 from account A to account B:
  1. **read(A)**
  2.  $A := A - 50$
  3. **write(A)**
  4. **read(B)**
  5.  $B := B + 50$
  6. **write(B)**
- Two main issues to deal with:
  - Failures of various kinds, such as hardware failures and system crashes
  - Concurrent execution of multiple transactions



# Example of Fund Transfer

- Transaction to transfer \$50 from account A to account B:
  1. **read(A)**
  2.  $A := A - 50$
  3. **write(A)**
  4. **read(B)**
  5.  $B := B + 50$
  6. **write(B)**
- **Atomicity requirement**
  - If the transaction fails after step 3 and before step 6, money will be “lost” leading to an inconsistent database state
    - Failure could be due to software or hardware
  - The system should ensure that updates of a partially executed transaction are not reflected in the database
- **Durability requirement** — once the user has been notified that the transaction has completed (i.e., the transfer of the \$50 has taken place), the updates to the database by the transaction must persist even if there are software or hardware failures.



# Example of Fund Transfer (Cont.)

- **Consistency requirement** in above example:
  - The sum of A and B is unchanged by the execution of the transaction
- In general, consistency requirements include
  - Explicitly specified integrity constraints such as primary keys and foreign keys
  - Implicit integrity constraints
    - e.g., sum of balances of all accounts, minus sum of loan amounts must equal value of cash-in-hand
  - A transaction must see a consistent database.
  - During transaction execution the database may be temporarily inconsistent.
  - When the transaction completes successfully the database must be consistent
    - Erroneous transaction logic can lead to inconsistency



# Example of Fund Transfer (Cont.)

- **Isolation requirement** — if between steps 3 and 6, another transaction T2 is allowed to access the partially updated database, it will see an inconsistent database (the sum  $A + B$  will be less than it should be).

T1

1. **read(A)**
2.  $A := A - 50$
3. **write(A)**
4. **read(B)**
5.  $B := B + 50$
6. **write(B)**

T2

read(A), read(B), print(A+B)

- Isolation can be ensured trivially by running transactions **serially**
  - That is, one after the other.
- However, executing multiple transactions concurrently has significant benefits, as we will see later.



# ACID Properties

A **transaction** is a unit of program execution that accesses and possibly updates various data items. To preserve the integrity of data the database system must ensure:

- **Atomicity.** Either all operations of the transaction are properly reflected in the database or none are.
- **Consistency.** Execution of a transaction in isolation preserves the consistency of the database.
- **Isolation.** Although multiple transactions may execute concurrently, each transaction must be unaware of other concurrently executing transactions. Intermediate transaction results must be hidden from other concurrently executed transactions.
  - That is, for every pair of transactions  $T_i$  and  $T_j$ , it appears to  $T_i$  that either  $T_j$  finished execution before  $T_i$  started, or  $T_j$  started execution after  $T_i$  finished.
- **Durability.** After a transaction completes successfully, the changes it has made to the database persist, even if there are system failures.



# Transaction State

- **Active** – the initial state; the transaction stays in this state while it is executing
- **Partially committed** – after the final statement has been executed.
- **Failed** -- after the discovery that normal execution can no longer proceed.
- **Aborted** – after the transaction has been rolled back and the database restored to its state prior to the start of the transaction. Two options after it has been aborted:
  - Restart the transaction
    - Can be done only if no internal logical error
  - Kill the transaction
- **Committed** – after successful completion.

# Atomicity Durability

# Atomicity

A transaction is a logical unit of work that must be either entirely completed or entirely undone. (All writes happen or none of them happen)

OK. What does this mean? Consider some pseudo-code

```
def transfer(source_acct_id, target_acct_id, amount)
```

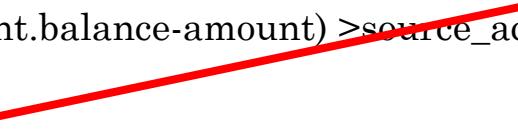
1. Check that both accounts exist.
2. IF *is\_checking\_account(source\_acct\_id)*
  1. Check that (source\_acct.balance-amount) > source\_account.overdraft\_limit
3. ELSE
  1. Check that (source\_count.balance-amount) >source\_account.minimum\_balance
4. Update source account
5. Update target account.
6. INSERT a record into transfer tracking table.

# Atomicity

A transaction is a logical unit of work that must be either entirely completed or entirely undone. (All writes happen or none of them happen)

OK. What does this mean? Consider some pseudo-code

```
def transfer(source_acct_id, target_acct_id, amount)
```

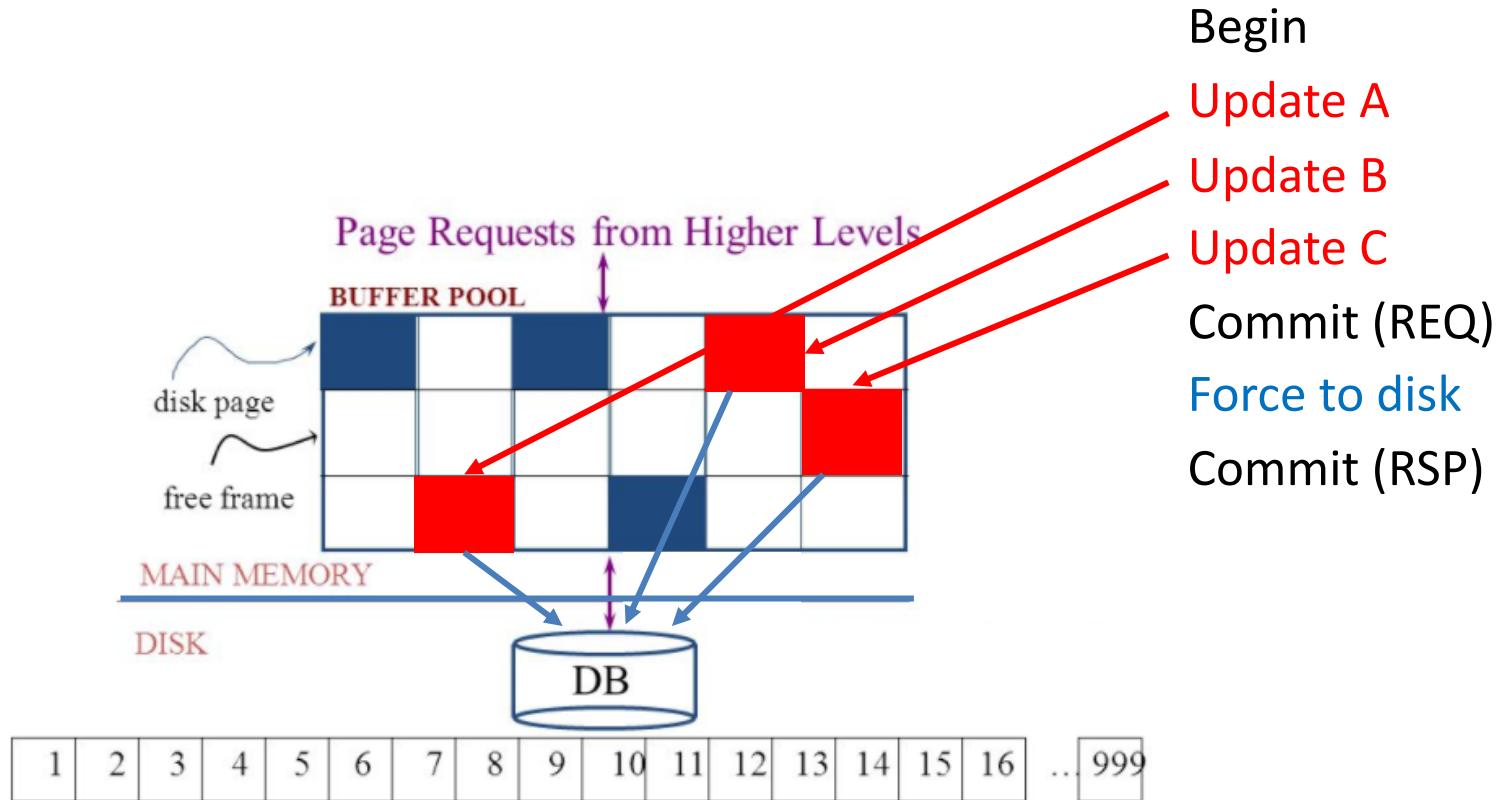
1. Check that both accounts exist.
2. IF *is\_checking\_account(source\_acct\_id)*
  1. Check that (source\_acct.balance-amount) > source\_acct.balance
3. ELSE
  1. Check that (source\_count.balance-amount) >source\_acct.balance
4. Update source account
5. Update target account. 
6. INSERT a record into transfer tracking table.



# Atomicty

- Transaction programs and databases are fast (milliseconds).  
What are the chances of the failure occurring in the wrong spot?
- Well, that doesn't really matter. If it happens,
  - Someone lost money and
  - There is no record off it. Someone is going to very upset.
- Even a small server can have thousands of concurrent transactions →
  - There will be corruptions because some transaction will be in the wrong place at the wrong time.
  - Unless we do something in the DBMS
  - Because HW and software inevitably fail
  - And sadly, SW is especially prone to failure when under load

# Simplistic Approach



# Simplistic Approach

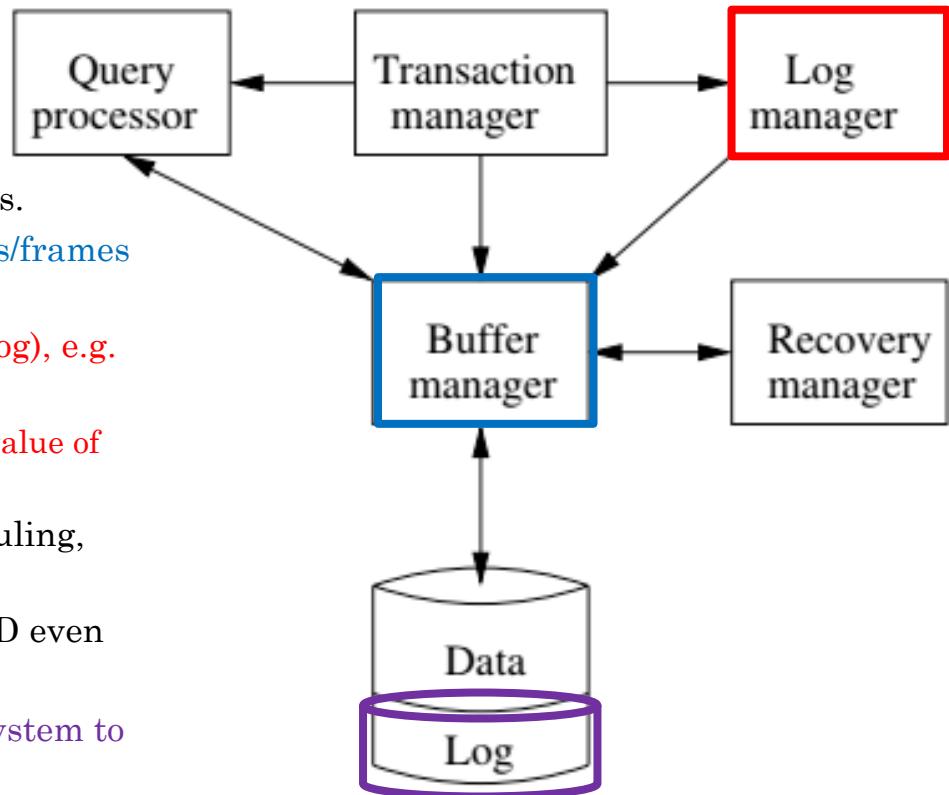
There are several problems with the simplistic approach.

1. The approach does not solve the problem
  1. Some writes might succeed.
  2. Some might be interrupted by the failure, or require retry.
2. Writes may be random and scattered. N updates might
  1. Change a few bytes in N data frames
  2. A few bytes in M index framesTransaction rate becomes bottlenecked by write I/O rate, even though a relative small number of bytes change/transaction.
3. Written frames must be held in memory.
  1. Lots of transactions
  2. Randomly writing small pieces of lots of frames.
  3. Consumes lots of memory with pinned pages.
  4. Degrades the performance and optimization of the buffer.
    1. The optimal buffer replacement policy wants to hold frames that will be reused.
    2. Not frames that have been touched and never reused.

# DBMS ACID Implementation

## Implementation Subsystems

- *Query processor* schedules and executes queries.
- *Buffer manager* controls reading/writing blocks/frames to/from disk.
- *Log manager* journals/records events to disk (log), e.g.
  - Transaction start/commit/abort
  - Transaction update of a block and previous value of data.
- *Transaction manager* coordinates query scheduling, buffer read/write and logging to ensure ACID.
- *Recovery manager* processes log to ensure ACID even after transaction or system failures.
- *Log* is a special type of block file used by the system to optimize performance and ensure ACID.

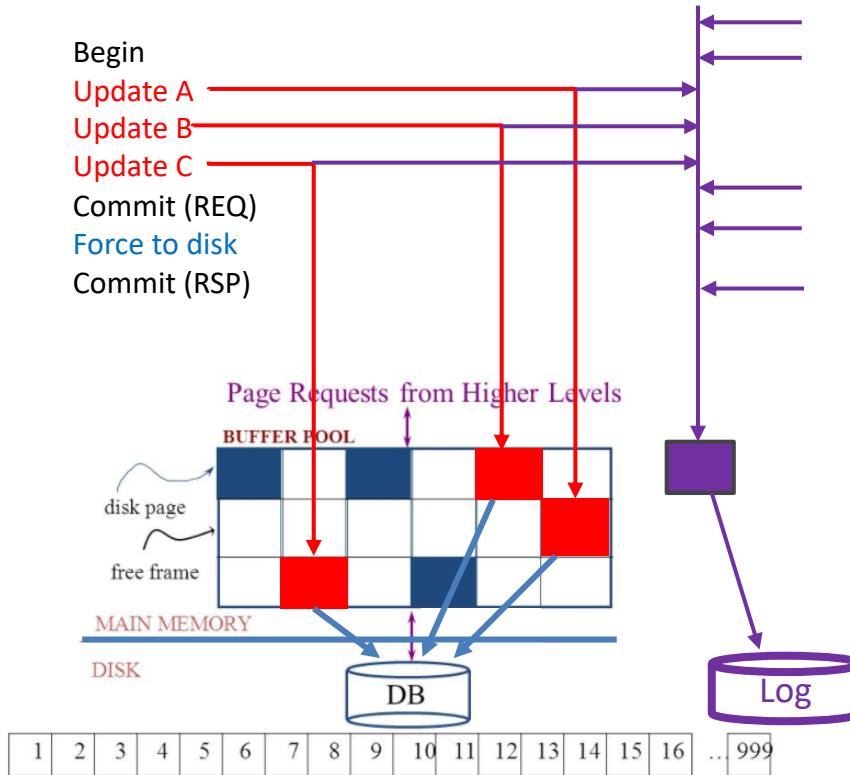


# Logging

The DBMS logs every transaction event

- *Log Sequence Number (LSN)*: A unique ID for a log record.
- *Prev LSN*: A link to their last log record.
- *Transaction ID number*.
- *Type*: Describes the type of database log record.
  - **Update Log Record**
    - *PageID*: A reference to the Page ID of the modified page.
    - *Length and Offset*: Length in bytes and offset of the page are usually included.
    - *Before and After Images* of records.
  - **Compensation Log Record**
  - **Commit Record**
  - **Abort Record Checkpoint Record**
  - **Completion Record** notes that all work has been done for this particular transaction.

# Write Ahead Logging



## DBMS (Redo processing)

- Write log events from all transactions into a single log stream.
- Multiple events per page
- Forces (writes) log record on COMMIT/ABORT
  - Single block I/O records many updates
  - Versus multiple block I/Os, each recording a single change.
  - All of a transaction's updates recorded in one I/O versus many.
- If there is a failure
  - DBMS sequentially reads log.
  - Applies changes to modified pages that were not saved to disk.
  - Then resumes normal processing.

# Write Ahead Logging

- Force every write to disk?
  - Poor response time.
  - But provides durability.
- Steal buffer-pool frames from uncommitted transactions?
  - If not, poor performance/caching performance
  - If yes, how can we ensure atomicity?  
Uncommitted updates on disk

	No Steal	Steal
Force	Trivial	
No Force		Desired

## DBMS (Undo processing)

- Enable steal policy to improve cache performance by
  - Avoiding lots of pinned pages
  - Unlikely to be reused soon.
- Before stealing
  - Force log record to disk.
  - Update log entry has data record
    - Before image
    - After image
- If there is a failure
  - DBMS sequentially reads log.
  - Undoes changes to
    - modified pages, uncommitted pages
    - That were saved to disk.
  - Then resumes normal processing.

ARIES recovery involves three passes

## 1. Analysis pass:

- Determine which transactions to undo
- Determine which pages were dirty (disk version not up to date) at time of
- RedoLSN: LSN from which redo should start

## 2. Redo pass:

- Repeats history, redoing all actions from RedoLSN  
(updated committed but not written changes to pages)
- RecLSN and PageLSNs are used to avoid redoing actions already reflected on page

## 3. Undo pass:

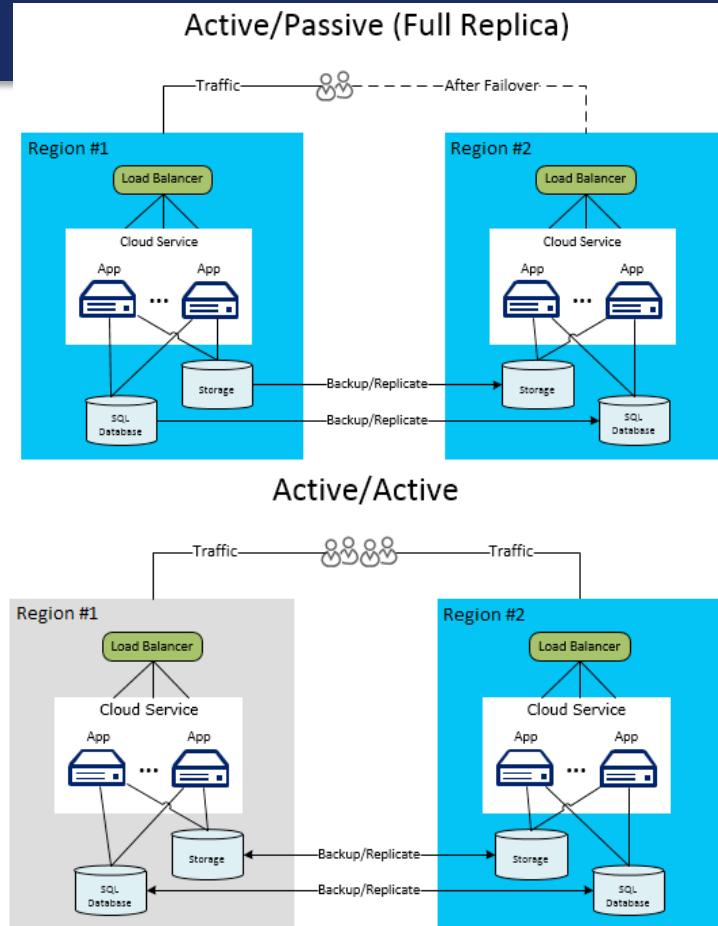
- Rolls back all incomplete transactions (with uncommitted pages written to disk).
- Transactions whose abort was complete earlier are not undone

# Durability

- Write changes to disk trivially achieves durability.
- DBMS engine uses write-ahead-logging to
  - Achieve durability
  - But with better performance through more efficient caching and I/O.
- Well, disks fail. How is that durable.
  - RAID and other solutions.
  - Disk subsystems, including entire RAID device, fail →
    - Duplex writes
    - To independent disk subsystems.
- Well, there are earthquakes, floods, etc.

# Availability and Replication

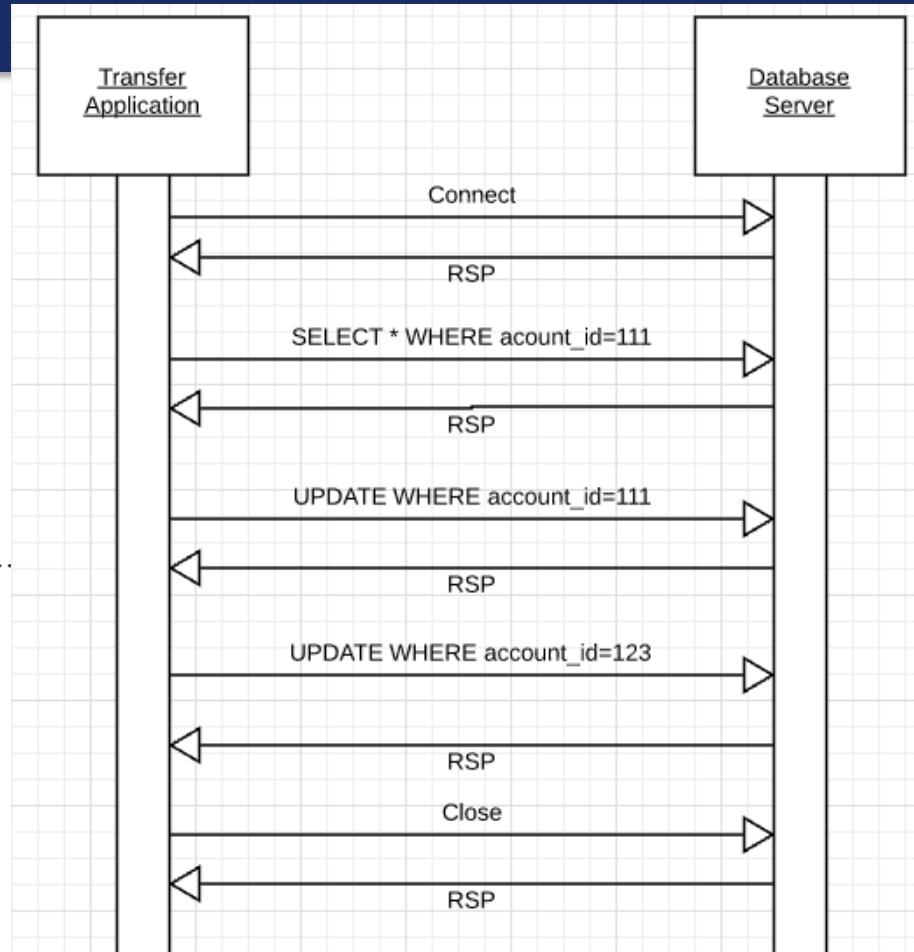
- There are two basic patterns
  - Active/Passive
    - All requests go to *master* during normal processing.
    - Updates are transactionally queued for processing at passive backup.
    - Failure of *master*
      - Routes subsequent requests to *backup*.
      - Backup must process and commit updates before accepting requests.
  - Active/Active
    - Both environments process requests.
    - Some form of distributed transaction commit required to synchronize updates on both copies.
- Multi-system communication to guarantee consistency is the foundation for tradeoffs in CAP.
  - The system can be CAP if and only iff
  - There are never any partitions or system failures
  - Which is unrealistic in cloud/Internet systems.



# *Isolation*

# Isolation

- Transfer \$50 from
  - account\_id=111 to
  - account\_id=123
- Requires 3 SQL statements
  - SELECT from 111 to check balance  $\geq \$50$
  - UPDATE account\_id=111
  - UPDATE account\_id=123
- There are some interesting scenarios
  - Two different programs read the balance (\$51)
  - And decide removing \$50 is OK.
- DB constraints can prevent the conflict from happening, but ...
  - There are more complex scenarios that constraints do not prevent.
  - Not ALL databases support constraints.
  - The “correct” execution should be that
    - One transaction responds “insufficient funds”
    - Before attempting transfer instead of after attempting.



# Isolation

- Try to transfer \$100 from account A to account B
  - Consider two simultaneous transfer transactions T1 and T2.
  - There are two equally **correct** executions
    1. T1 transfers, T2 responds “insufficient funds” and does not attempt transfer
    2. T2 transfers, T1 responds “insufficient funds” and does not attempt transfer
  - Each correct simultaneous execution is equivalent to a serial (sequential) execution schedule
    - (1) Execute T1, Execute T2
    - (2) Execute T2, Execute T1
  - NOTE:
    - We are focusing on correctness not
    - Fairness:
      - We do not care which transaction was actually submitted first.
      - And probably do not know due to networking, etc.

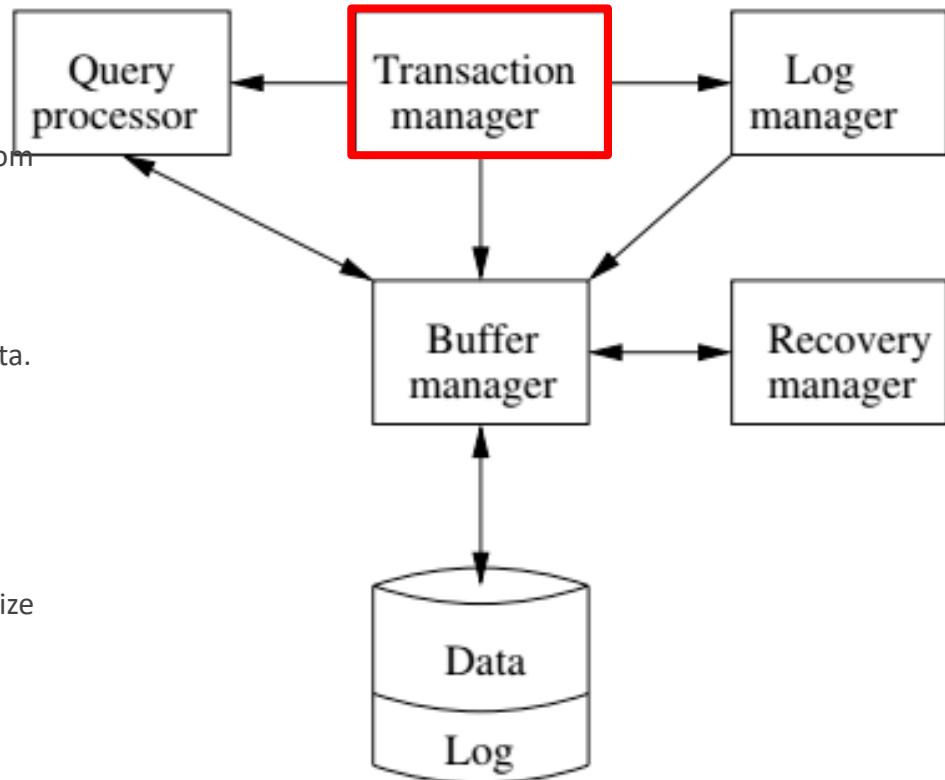
# Serializability

“In [concurrency control](#) of [databases](#),<sup>[1][2]</sup> [transaction processing](#) (transaction management), and various [transactional](#) applications (e.g., [transactional memory](#)<sup>[3]</sup> and [software transactional memory](#)), both centralized and [distributed](#), a transaction [schedule](#) is **serializable** if its outcome (e.g., the resulting database state) is equal to the outcome of its transactions executed serially, i.e. without overlapping in time. Transactions are normally executed concurrently (they overlap), since this is the most efficient way. Serializability is the major correctness criterion for concurrent transactions' executions. It is considered the highest level of [isolation](#) between [transactions](#), and plays an essential role in [concurrency control](#). As such it is supported in all general purpose database systems.”  
(<https://en.wikipedia.org/wiki/Serializability>)

# DBMS ACID Implementation

## Implementation Subsystems

- *Query processor* schedules and executes queries.
- *Buffer manager* controls reading/writing blocks/frames to/from disk.
- *Log manager* journals/records events to disk (log), e.g.
  - Transaction start/commit/abort
  - Transaction update of a block and previous value of data.
- *Transaction manager* coordinates query scheduling, buffer read/write and logging to ensure ACID.
- *Recovery manager* processes log to ensure ACID even after transaction or system failures.
- *Log* is a special type of block file used by the system to optimize performance and ensure ACID.



Garcia-Molina et al., p. 846

# Schedule

## 18.1.1 Schedules

A *schedule* is a sequence of the important actions taken by one or more transactions. When studying concurrency control, the important read and write actions take place in the main-memory buffers, not the disk. That is, a database element  $A$  that is brought to a buffer by some transaction  $T$  may be read or written in that buffer not only by  $T$  but by other transactions that access  $A$ .

$T_1$	$T_2$
READ(A,t)	READ(A,s)
$t := t+100$	$s := s*2$
WRITE(A,t)	WRITE(A,s)
READ(B,t)	READ(B,s)
$t := t+100$	$s := s*2$
WRITE(B,t)	WRITE(B,s)

Figure 18.2: Two transactions

Garcia-Molina et al.

## 18.1.2 Serial Schedules

- Assume there are three

- concurrently executing transactions allowed.
- T<sub>1</sub>, T<sub>2</sub> and T<sub>3</sub>

- The transaction manager

- Enables concurrent execution
- But schedules individual operations
- To ensure that the final DB state
- Is *equivalent* to one of the following schedules
  - T<sub>1</sub>, T<sub>2</sub>, T<sub>3</sub>
  - T<sub>1</sub>, T<sub>3</sub>, T<sub>2</sub>
  - T<sub>2</sub>, T<sub>1</sub>, T<sub>3</sub>
  - T<sub>2</sub>, T<sub>3</sub>, T<sub>1</sub>
  - T<sub>3</sub>, T<sub>1</sub>, T<sub>2</sub>
  - T<sub>3</sub>, T<sub>2</sub>, T<sub>1</sub>

A schedule is *serial* if its actions consist of all the actions of one transaction, then all the actions of another transaction, and so on. No mixing of the actions

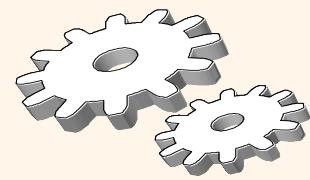
	T <sub>1</sub>	T <sub>2</sub>	A	B
	READ(A,t)		25	25
	t := t+100			
	WRITE(A,t)		125	
	READ(B,t)			
	t := t+100			
	WRITE(B,t)		125	
		READ(A,s)		
		s := s*2		
		WRITE(A,s)	250	
		READ(B,s)		
		s := s*2		
		WRITE(B,s)	250	

Concurrent execution was *serializable*.

Figure 18.3: Serial schedule in which T<sub>1</sub> precedes T<sub>2</sub>

# Serializability ([en.wikipedia.org/wiki/Serializability](https://en.wikipedia.org/wiki/Serializability))

- **Serializability** is used to keep the data in the data item in a consistent state. Serializability is a property of a transaction schedule (history). It relates to the isolation property of a database transaction.
- **Serializability** of a schedule means equivalence (in the outcome, the database state, data values) to a *serial schedule* (i.e., sequential with no transaction overlap in time) with the same transactions. It is the major criterion for the correctness of concurrent transactions' schedule, and thus supported in all general purpose database systems.
- **The rationale behind serializability** is the following:
  - If each transaction is correct by itself, i.e., meets certain integrity conditions,
  - then a schedule that comprises any *serial* execution of these transactions is correct (its transactions still meet their conditions):
    - "Serial" means that transactions do not overlap in time and cannot interfere with each other, i.e., complete *isolation* between each other exists.
    - Any order of the transactions is legitimate, (...)
    - As a result, a schedule that comprises any execution (not necessarily serial) that is equivalent (in its outcome) to any serial execution of these transactions, is correct.



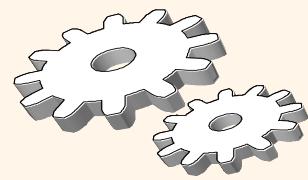
# Lock-Based Concurrency Control

## ❖ Strict Two-phase Locking (Strict 2PL) Protocol:

- Each Xact must obtain a **S (shared) lock** on object before reading, and an **X (exclusive) lock** on object before writing.
- All locks held by a transaction are released when the transaction completes
  - **(Non-strict) 2PL Variant:** Release locks anytime, but cannot acquire locks after releasing any lock.
- If an Xact holds an X lock on an object, no other Xact can get a lock (S or X) on that object.

## ❖ Strict 2PL allows only serializable schedules.

- Additionally, it simplifies transaction aborts
- **(Non-strict) 2PL** also allows only serializable schedules, but involves more complex abort processing



# Aborting a Transaction

- ❖ If a transaction  $T_i$  is aborted, all its actions have to be undone. Not only that, if  $T_j$  reads an object last written by  $T_i$ ,  $T_j$  must be aborted as well!
- ❖ Most systems avoid such *cascading aborts* by releasing a transaction's locks only at commit time.
  - If  $T_i$  writes an object,  $T_j$  can read this only after  $T_i$  commits.
- ❖ In order to *undo* the actions of an aborted transaction, the DBMS maintains a *log* in which every write is recorded. This mechanism is also used to recover from system crashes: all active Xacts at the time of the crash are aborted when the system comes back up.

# MySQL (Locking) Isolation

## 13.3.6 SET TRANSACTION Syntax

```
1  SET [GLOBAL | SESSION] TRANSACTION  
2      transaction_characteristic [, transaction_characteristic] ...  
3  
4  transaction_characteristic:  
5      ISOLATION LEVEL level  
6      | READ WRITE  
7      | READ ONLY  
8  
9  level:  
10     REPEATABLE READ  
11     | READ COMMITTED  
12     | READ UNCOMMITTED  
13     | SERIALIZABLE
```

### Scope of Transaction Characteristics

You can set transaction characteristics globally, for the current session, or for the next transaction:

- With the `GLOBAL` keyword, the statement applies globally for all subsequent sessions. Existing sessions are unaffected.
- With the `SESSION` keyword, the statement applies to all subsequent transactions performed within the current session.
- Without any `SESSION` or `GLOBAL` keyword, the statement applies to the next (not started) transaction performed within the current session. Subsequent transactions revert to using the `SESSION` isolation level.

# Isolation Levels

([https://en.wikipedia.org/wiki/Isolation\\_\(database\\_systems\)](https://en.wikipedia.org/wiki/Isolation_(database_systems)))

Not all transaction use cases require 2PL and serializable execution. Databases support a set of levels.

- **Serializable**
  - With a lock-based [concurrency control](#) DBMS implementation, [serializability](#) requires read and write locks (acquired on selected data) to be released at the end of the transaction. Also *range-locks* must be acquired when a [SELECT](#) query uses a ranged *WHERE* clause, especially to avoid the [phantom reads](#) phenomenon.
  - *The execution of concurrent SQL-transactions at isolation level SERIALIZABLE is guaranteed to be serializable. A serializable execution is defined to be an execution of the operations of concurrently executing SQL-transactions that produces the same effect as some serial execution of those same SQL-transactions. A serial execution is one in which each SQL-transaction executes to completion before the next SQL-transaction begins.*
- **Repeatable reads**
  - In this isolation level, a lock-based [concurrency control](#) DBMS implementation keeps read and write locks (acquired on selected data) until the end of the transaction. However, *range-locks* are not managed, so [phantom reads](#) can occur.
  - Write skew is possible at this isolation level, a phenomenon where two writes are allowed to the same column(s) in a table by two different writers (who have previously read the columns they are updating), resulting in the column having data that is a mix of the two transactions.<sup>[3][4]</sup>
- **Read committed**
  - In this isolation level, a lock-based [concurrency control](#) DBMS implementation keeps write locks (acquired on selected data) until the end of the transaction, but read locks are released as soon as the [SELECT](#) operation is performed (so the [non-repeatable reads phenomenon](#) can occur in this isolation level). As in the previous level, *range-locks* are not managed.
  - Putting it in simpler words, read committed is an isolation level that guarantees that any data read is committed at the moment it is read. It simply restricts the reader from seeing any intermediate, uncommitted, 'dirty' read. It makes no promise whatsoever that if the transaction re-issues the read, it will find the same data; data is free to change after it is read.
- **Read uncommitted**
  - This is the *lowest* isolation level. In this level, [dirty reads](#) are allowed, so one transaction may see *not-yet-committed* changes made by other transactions

# In Databases, Cursors Define *Isolation*

- We have talked about ACID transactions

Isolation level	Dirty reads	Non-repeatable reads	Phantoms
Read Uncommitted	may occur	may occur	may occur
Read Committed	-	may occur	may occur
Repeatable Read	-	-	may occur
Serializable	-	-	-

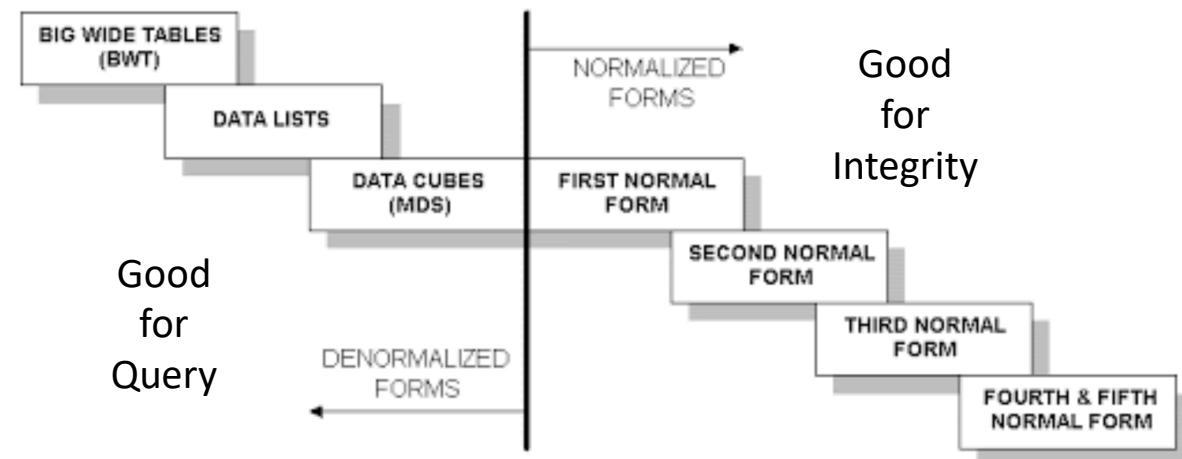
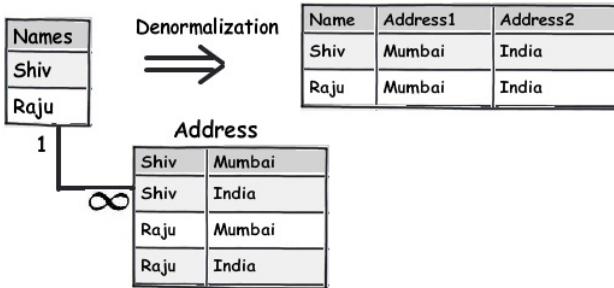
- Isolation
  - Determines what happens when two or more threads are manipulating the data at the same time.
  - And is defined relative to where cursors are and what they have touched.
  - Because the cursor movement determines *what you are reading or have read*.
- *But, ... Cursors are client conversation state and cannot be used in REST.*

```
InitialContext ctx = new InitialContext();
DataSource ds = (DataSource)
ctx.lookup("jdbc/MyBase");
Connection con = ds.getConnection();
DatabaseMetaData dbmd = con.getMetaData();
if (dbmd.supportsTransactionIsolationLevel(TRANSACTION_SERIALIZABLE))
{ Connection.setTransactionIsolation(TRANSACTION_SERIALIZABLE); }
```

# Decision Support

*Isolation*

# Wide Flat Tables



- Improve query performance by precomputing and saving:
  - JOINs
  - Aggregation
  - Derived/computed columns
- One of the primary strength of the relational model is maintaining “integrity” when applications create, update and delete data. This relies on:
  - The core capabilities of the relational model, e.g. constraints.
  - A well-design database (We will cover a formal definition – “normalization” in more detail later.)
- Data models that are well designed for integrity are very bad for read only analysis queries.  
We will build and analyze wide flat tables as part of the analysis tasks in HW3, HW4 as projects.



# Chapter 20: Data Analysis

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Chapter 20: Data Analysis

- Decision Support Systems
- Data Warehousing
- Data Mining
- Classification
- Association Rules
- Clustering



# Decision Support Systems

- **Decision-support systems** are used to make business decisions, often based on data collected by on-line transaction-processing systems.
- Examples of business decisions:
  - What items to stock?
  - What insurance premium to change?
  - To whom to send advertisements?
- Examples of data used for making decisions
  - Retail sales transaction details
  - Customer profiles (income, age, gender, etc.)



# Decision-Support Systems: Overview

- **Data analysis** tasks are simplified by specialized tools and SQL extensions
  - Example tasks
    - ▶ For each product category and each region, what were the total sales in the last quarter and how do they compare with the same quarter last year
    - ▶ As above, for each product category and each customer category
- **Statistical analysis** packages (e.g., : S++) can be interfaced with databases
  - Statistical analysis is a large field, but not covered here
- **Data mining** seeks to discover knowledge automatically in the form of statistical rules and patterns from large databases

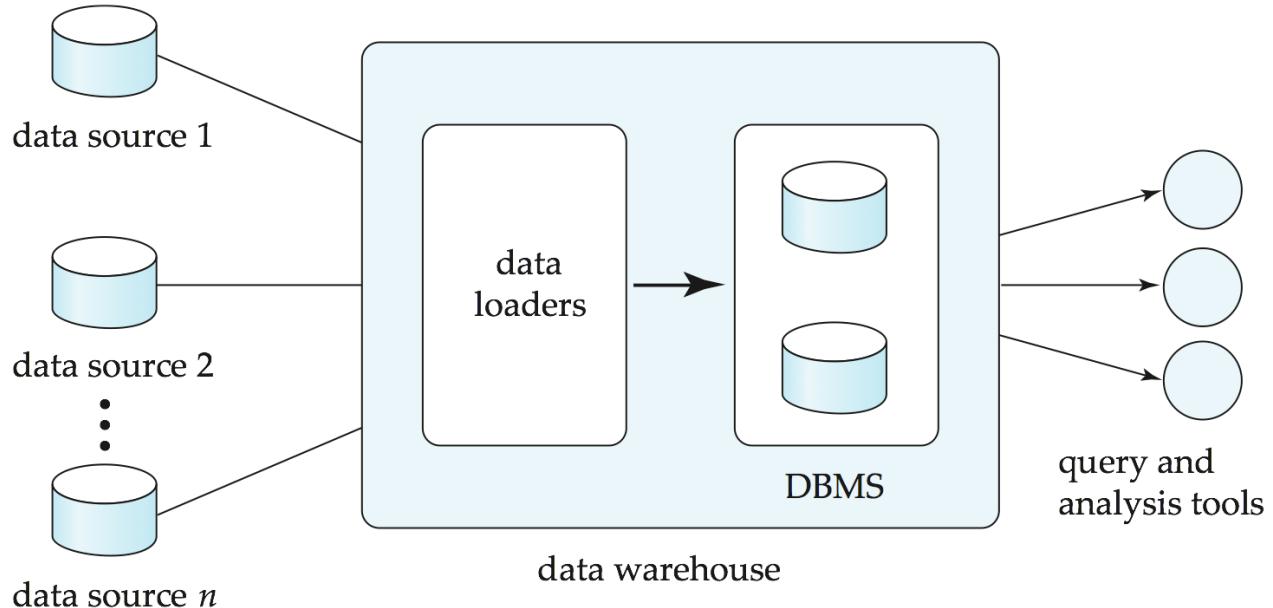


# Data Warehousing

- Data sources often store only current data, not historical data
- Corporate decision making requires a unified view of all organizational data, including historical data
- A **data warehouse** is a repository (archive) of information gathered from multiple sources, stored under a unified schema, at a single site
  - Greatly simplifies querying, permits study of historical trends
  - Shifts decision support query load away from transaction processing systems



# Data Warehousing





# Design Issues

## ■ When and how to gather data

- **Source driven architecture**: data sources transmit new information to warehouse, either continuously or periodically (e.g., at night)
- **Destination driven architecture**: warehouse periodically requests new information from data sources
- Keeping warehouse exactly synchronized with data sources (e.g., using two-phase commit) is too expensive
  - ▶ Usually OK to have slightly out-of-date data at warehouse
  - ▶ Data/updates are periodically downloaded from online transaction processing (OLTP) systems.

## ■ What schema to use

- Schema integration



# More Warehouse Design Issues

## ■ *Data cleansing*

- E.g., correct mistakes in addresses (misspellings, zip code errors)
- **Merge** address lists from different sources and **purge** duplicates

## ■ *How to propagate updates*

- Warehouse schema may be a (materialized) view of schema from data sources

## ■ *What data to summarize*

- Raw data may be too large to store on-line
- Aggregate values (totals/subtotals) often suffice
- Queries on raw data can often be transformed by query optimizer to use aggregate values

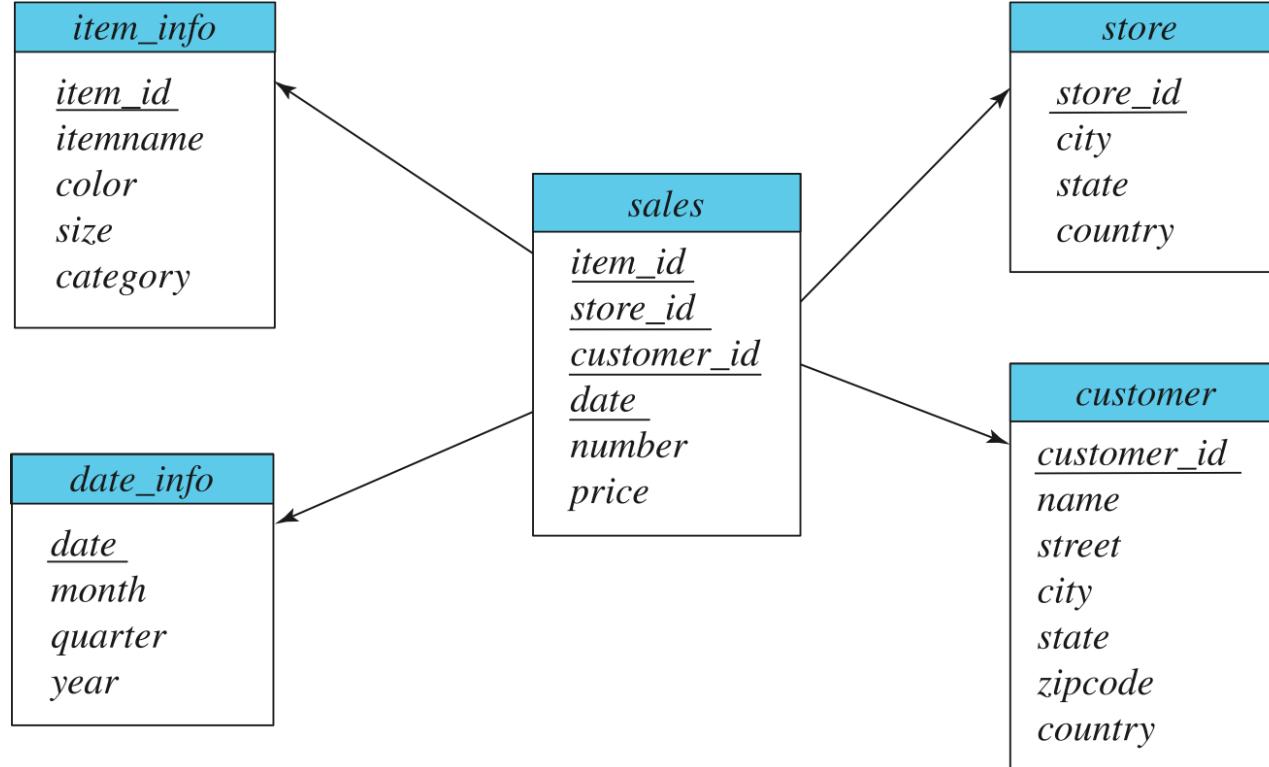


# Warehouse Schemas

- Dimension values are usually encoded using small integers and mapped to full values via dimension tables
- Resultant schema is called a **star schema**
  - More complicated schema structures
    - ▶ **Snowflake schema**: multiple levels of dimension tables
    - ▶ **Constellation**: multiple fact tables



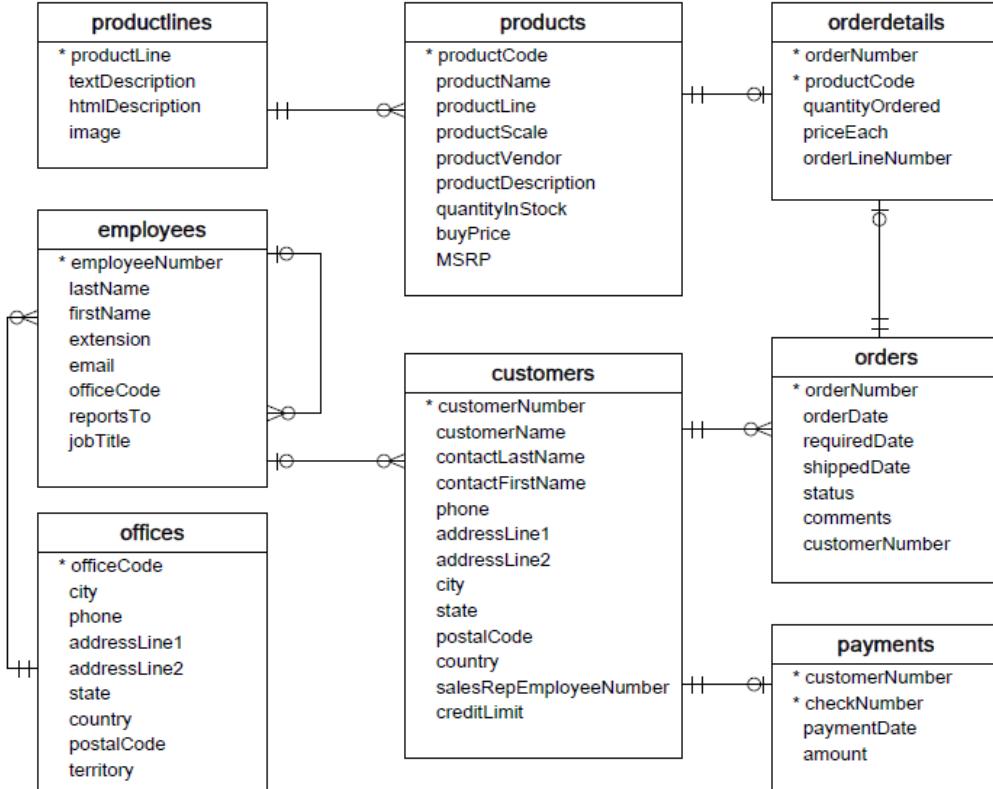
# Data Warehouse Schema



# *Worked Example*

# *Classic Models*

# Classic Models



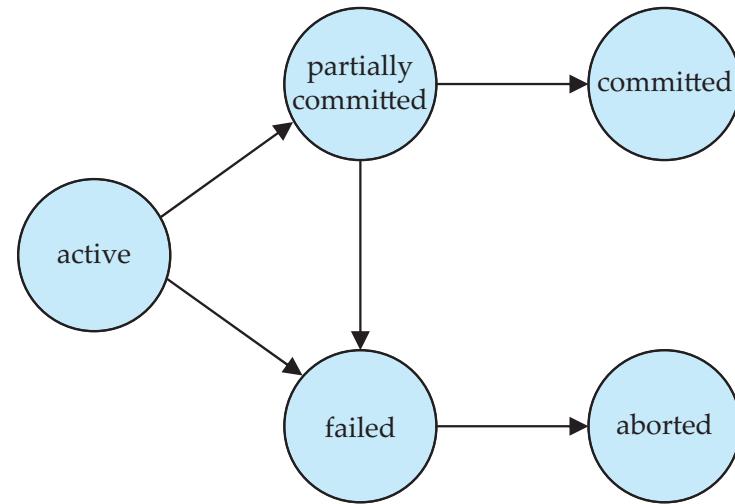
- This schema is “good” relative
  - To normalization.
  - Avoiding update anomalies.
- Common, intuitive decision support and insight queries require a lot of
  - JOINs
  - Aggregates
  - ... ...
  - These operations are slow and can interfere with operational transactions.
- We will look at some transformation to a decision support schema.

# Go To Notebook

# Transaction Details (Skip, but students to read)



# Transaction State (Cont.)





# Concurrent Executions

- Multiple transactions are allowed to run concurrently in the system.  
Advantages are:
  - **Increased processor and disk utilization**, leading to better transaction *throughput*
    - E.g., one transaction can be using the CPU while another is reading from or writing to the disk
  - **Reduced average response time** for transactions: short transactions need not wait behind long ones.
- **Concurrency control schemes** – mechanisms to achieve isolation
  - That is, to control the interaction among the concurrent transactions in order to prevent them from destroying the consistency of the database
    - Will study in Chapter 15, after studying notion of correctness of concurrent executions.



# Schedules

- **Schedule** – a sequences of instructions that specify the chronological order in which instructions of concurrent transactions are executed
  - A schedule for a set of transactions must consist of all instructions of those transactions
  - Must preserve the order in which the instructions appear in each individual transaction.
- A transaction that successfully completes its execution will have a commit instruction as the last statement
  - By default transaction assumed to execute commit instruction as its last step
- A transaction that fails to successfully complete its execution will have an abort instruction as the last statement



# Schedule 1

- Let  $T_1$  transfer \$50 from  $A$  to  $B$ , and  $T_2$  transfer 10% of the balance from  $A$  to  $B$ .
- A **serial** schedule in which  $T_1$  is followed by  $T_2$ :

$T_1$	$T_2$
read ( $A$ ) $A := A - 50$ write ( $A$ ) read ( $B$ ) $B := B + 50$ write ( $B$ ) commit	read ( $A$ ) $temp := A * 0.1$ $A := A - temp$ write ( $A$ ) read ( $B$ ) $B := B + temp$ write ( $B$ ) commit



## Schedule 2

- A serial schedule where  $T_2$  is followed by  $T_1$

$T_1$	$T_2$
read ( $A$ ) $A := A - 50$ write ( $A$ ) read ( $B$ ) $B := B + 50$ write ( $B$ ) commit	read ( $A$ ) $temp := A * 0.1$ $A := A - temp$ write ( $A$ ) read ( $B$ ) $B := B + temp$ write ( $B$ ) commit



## Schedule 3

- Let  $T_1$  and  $T_2$  be the transactions defined previously. The following schedule is not a serial schedule, but it is *equivalent* to Schedule 1

$T_1$	$T_2$
read ( $A$ ) $A := A - 50$ write ( $A$ )	read ( $A$ ) $temp := A * 0.1$ $A := A - temp$ write ( $A$ )
read ( $B$ ) $B := B + 50$ write ( $B$ ) commit	read ( $B$ ) $B := B + temp$ write ( $B$ ) commit

- In Schedules 1, 2 and 3, the sum  $A + B$  is preserved.



## Schedule 4

- The following concurrent schedule does not preserve the value of  $(A + B)$ .

$T_1$	$T_2$
read ( $A$ ) $A := A - 50$  write ( $A$ ) read ( $B$ ) $B := B + 50$ write ( $B$ ) commit	read ( $A$ ) $temp := A * 0.1$ $A := A - temp$ write ( $A$ ) read ( $B$ )  $B := B + temp$ write ( $B$ ) commit



# Serializability

- **Basic Assumption** – Each transaction preserves database consistency.
- Thus, serial execution of a set of transactions preserves database consistency.
- A (possibly concurrent) schedule is serializable if it is equivalent to a serial schedule. Different forms of schedule equivalence give rise to the notions of:
  1. **Conflict serializability**
  2. **View serializability**



## *Simplified view of transactions*

- We ignore operations other than **read** and **write** instructions
- We assume that transactions may perform arbitrary computations on data in local buffers in between reads and writes.
- Our simplified schedules consist of only **read** and **write** instructions.



# Conflicting Instructions

- Instructions  $I_i$  and  $I_j$  of transactions  $T_i$  and  $T_j$  respectively, **conflict** if and only if there exists some item  $Q$  accessed by both  $I_i$  and  $I_j$ , and at least one of these instructions wrote  $Q$ .
  1.  $I_i = \text{read}(Q)$ ,  $I_j = \text{read}(Q)$ .  $I_i$  and  $I_j$  don't conflict.
  2.  $I_i = \text{read}(Q)$ ,  $I_j = \text{write}(Q)$ . They conflict.
  3.  $I_i = \text{write}(Q)$ ,  $I_j = \text{read}(Q)$ . They conflict
  4.  $I_i = \text{write}(Q)$ ,  $I_j = \text{write}(Q)$ . They conflict
- Intuitively, a conflict between  $I_i$  and  $I_j$  forces a (logical) temporal order between them.
- If  $I_i$  and  $I_j$  are consecutive in a schedule and they do not conflict, their results would remain the same even if they had been interchanged in the schedule.



# Conflict Serializability

- If a schedule  $S$  can be transformed into a schedule  $S'$  by a series of swaps of non-conflicting instructions, we say that  $S$  and  $S'$  are **conflict equivalent**.
- We say that a schedule  $S$  is **conflict serializable** if it is conflict equivalent to a serial schedule



## Conflict Serializability (Cont.)

- Schedule 3 can be transformed into Schedule 6, a serial schedule where  $T_2$  follows  $T_1$ , by series of swaps of non-conflicting instructions. Therefore Schedule 3 is conflict serializable.

$T_1$	$T_2$
read ( $A$ ) write ( $A$ )	read ( $A$ ) write ( $A$ )
read ( $B$ ) write ( $B$ )	read ( $B$ ) write ( $B$ )

Schedule 3

$T_1$	$T_2$
read ( $A$ ) write ( $A$ ) read ( $B$ ) write ( $B$ )	read ( $B$ ) write ( $B$ )
	read ( $A$ ) write ( $A$ ) read ( $B$ ) write ( $B$ )

Schedule 6



## Conflict Serializability (Cont.)

- Example of a schedule that is not conflict serializable:

$T_3$	$T_4$
read ( $Q$ )	
write ( $Q$ )	write ( $Q$ )

- We are unable to swap instructions in the above schedule to obtain either the serial schedule  $\langle T_3, T_4 \rangle$ , or the serial schedule  $\langle T_4, T_3 \rangle$ .



# View Serializability

- Let  $S$  and  $S'$  be two schedules with the same set of transactions.  $S$  and  $S'$  are **view equivalent** if the following three conditions are met, for each data item  $Q$ ,
  1. If in schedule  $S$ , transaction  $T_i$  reads the initial value of  $Q$ , then in schedule  $S'$  also transaction  $T_i$  must read the initial value of  $Q$ .
  2. If in schedule  $S$  transaction  $T_i$  executes **read**( $Q$ ), and that value was produced by transaction  $T_j$  (if any), then in schedule  $S'$  also transaction  $T_i$  must read the value of  $Q$  that was produced by the same **write**( $Q$ ) operation of transaction  $T_j$ .
  3. The transaction (if any) that performs the final **write**( $Q$ ) operation in schedule  $S$  must also perform the final **write**( $Q$ ) operation in schedule  $S'$ .
- As can be seen, view equivalence is also based purely on **reads** and **writes** alone.



## View Serializability (Cont.)

- A schedule  $S$  is **view serializable** if it is view equivalent to a serial schedule.
- Every conflict serializable schedule is also view serializable.
- Below is a schedule which is view-serializable but *not* conflict serializable.

$T_{27}$	$T_{28}$	$T_{29}$
read ( $Q$ )		
write ( $Q$ )	write ( $Q$ )	write ( $Q$ )

- What serial schedule is above equivalent to?
- Every view serializable schedule that is not conflict serializable has **blind writes**.



# Other Notions of Serializability

- The schedule below produces same outcome as the serial schedule  $\langle T_1, T_5 \rangle$ , yet is not conflict equivalent or view equivalent to it.

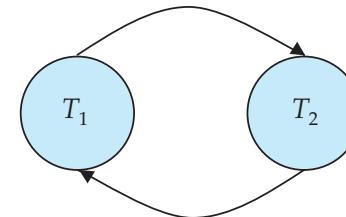
$T_1$	$T_5$
read ( $A$ )	
$A := A - 50$	
write ( $A$ )	
	read ( $B$ )
	$B := B - 10$
	write ( $B$ )
read ( $B$ )	
$B := B + 50$	
write ( $B$ )	
	read ( $A$ )
	$A := A + 10$
	write ( $A$ )

- Determining such equivalence requires analysis of operations other than read and write.



# Testing for Serializability

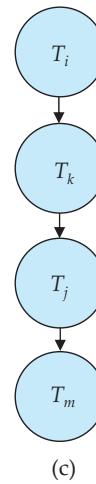
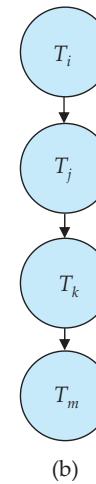
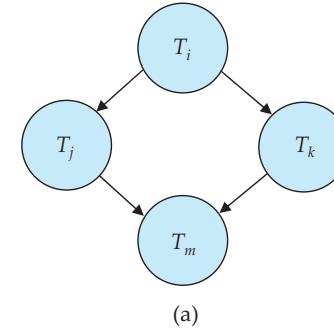
- Consider some schedule of a set of transactions  $T_1, T_2, \dots, T_n$
- **Precedence graph** — a direct graph where the vertices are the transactions (names).
- We draw an arc from  $T_i$  to  $T_j$  if the two transaction conflict, and  $T_i$  accessed the data item on which the conflict arose earlier.
- We may label the arc by the item that was accessed.
- Example of a precedence graph





# Test for Conflict Serializability

- A schedule is conflict serializable if and only if its precedence graph is acyclic.
- Cycle-detection algorithms exist which take order  $n^2$  time, where  $n$  is the number of vertices in the graph.
  - (Better algorithms take order  $n + e$  where  $e$  is the number of edges.)
- If precedence graph is acyclic, the serializability order can be obtained by a *topological sorting* of the graph.
  - This is a linear order consistent with the partial order of the graph.
  - For example, a serializability order for Schedule A would be  
 $T_5 \rightarrow T_1 \rightarrow T_3 \rightarrow T_2 \rightarrow T_4$ 
    - Are there others?





# Test for View Serializability

- The precedence graph test for conflict serializability cannot be used directly to test for view serializability.
  - Extension to test for view serializability has cost exponential in the size of the precedence graph.
- The problem of checking if a schedule is view serializable falls in the class of *NP*-complete problems.
  - Thus, existence of an efficient algorithm is *extremely* unlikely.
- However practical algorithms that just check some **sufficient conditions** for view serializability can still be used.



# Recoverable Schedules

Need to address the effect of transaction failures on concurrently running transactions.

- **Recoverable schedule** — if a transaction  $T_j$  reads a data item previously written by a transaction  $T_i$ , then the commit operation of  $T_i$  appears before the commit operation of  $T_j$ .
- The following schedule (Schedule 11) is not recoverable

$T_8$	$T_9$
read ( $A$ ) write ( $A$ )	
read ( $B$ )	read ( $A$ ) commit

- If  $T_8$  should abort,  $T_9$  would have read (and possibly shown to the user) an inconsistent database state. Hence, database must ensure that schedules are recoverable.



# Cascading Rollbacks

- **Cascading rollback** – a single transaction failure leads to a series of transaction rollbacks. Consider the following schedule where none of the transactions has yet committed (so the schedule is recoverable)

$T_{10}$	$T_{11}$	$T_{12}$
read ( $A$ ) read ( $B$ ) write ( $A$ )  abort	read ( $A$ ) write ( $A$ )	read ( $A$ )

If  $T_{10}$  fails,  $T_{11}$  and  $T_{12}$  must also be rolled back.

- Can lead to the undoing of a significant amount of work



# Cascadeless Schedules

- **Cascadeless schedules** — cascading rollbacks cannot occur;
  - For each pair of transactions  $T_i$  and  $T_j$  such that  $T_j$  reads a data item previously written by  $T_i$ , the commit operation of  $T_i$  appears before the read operation of  $T_j$ .
- Every Cascadeless schedule is also recoverable
- It is desirable to restrict the schedules to those that are cascadeless



# Concurrency Control

- A database must provide a mechanism that will ensure that all possible schedules are
  - either conflict or view serializable, and
  - are recoverable and preferably cascadeless
- A policy in which only one transaction can execute at a time generates serial schedules, but provides a poor degree of concurrency
  - Are serial schedules recoverable/cascadeless?
- Testing a schedule for serializability *after* it has executed is a little too late!
- **Goal** – to develop concurrency control protocols that will assure serializability.



# Concurrency Control (Cont.)

- Schedules must be conflict or view serializable, and recoverable, for the sake of database consistency, and preferably cascadeless.
- A policy in which only one transaction can execute at a time generates serial schedules, but provides a poor degree of concurrency.
- Concurrency-control schemes tradeoff between the amount of concurrency they allow and the amount of overhead that they incur.
- Some schemes allow only conflict-serializable schedules to be generated, while others allow view-serializable schedules that are not conflict-serializable.



# Concurrency Control vs. Serializability Tests

- Concurrency-control protocols allow concurrent schedules, but ensure that the schedules are conflict/view serializable, and are recoverable and cascadeless .
- Concurrency control protocols (generally) do not examine the precedence graph as it is being created
  - Instead a protocol imposes a discipline that avoids non-serializable schedules.
  - We study such protocols in Chapter 16.
- Different concurrency control protocols provide different tradeoffs between the amount of concurrency they allow and the amount of overhead that they incur.
- Tests for serializability help us understand why a concurrency control protocol is correct.



# Weak Levels of Consistency

- Some applications are willing to live with weak levels of consistency, allowing schedules that are not serializable
  - E.g., a read-only transaction that wants to get an approximate total balance of all accounts
  - E.g., database statistics computed for query optimization can be approximate (why?)
  - Such transactions need not be serializable with respect to other transactions
- Tradeoff accuracy for performance



# Levels of Consistency in SQL-92

- **Serializable** — default
- **Repeatable read** — only committed records to be read.
  - Repeated reads of same record must return same value.
  - However, a transaction may not be serializable – it may find some records inserted by a transaction but not find others.
- **Read committed** — only committed records can be read.
  - Successive reads of record may return different (but committed) values.
- **Read uncommitted** — even uncommitted records may be read.



# Levels of Consistency

- Lower degrees of consistency useful for gathering approximate information about the database
- Warning: some database systems do not ensure serializable schedules by default
- E.g., Oracle (and PostgreSQL prior to version 9) by default support a level of consistency called snapshot isolation (not part of the SQL standard)



# Transaction Definition in SQL

- In SQL, a transaction begins implicitly.
- A transaction in SQL ends by:
  - **Commit work** commits current transaction and begins a new one.
  - **Rollback work** causes current transaction to abort.
- In almost all database systems, by default, every SQL statement also commits implicitly if it executes successfully
  - Implicit commit can be turned off by a database directive
    - E.g., in JDBC -- `connection.setAutoCommit(false);`
- Isolation level can be set at database level
- Isolation level can be changed at start of transaction
  - E.g. In SQL **set transaction isolation level serializable**
  - E.g. in JDBC -- `connection.setTransactionIsolation(`  
`Connection.TRANSACTION_SERIALIZABLE)`



# Implementation of Isolation Levels

- Locking
  - Lock on whole database vs lock on items
  - How long to hold lock?
  - Shared vs exclusive locks
- Timestamps
  - Transaction timestamp assigned e.g. when a transaction begins
  - Data items store two timestamps
    - Read timestamp
    - Write timestamp
  - Timestamps are used to detect out of order accesses
- Multiple versions of each data item
  - Allow transactions to read from a “snapshot” of the database



# Transactions as SQL Statements

- E.g., Transaction 1:  
`select ID, name from instructor where salary > 90000`
- E.g., Transaction 2:  
`insert into instructor values ('11111', 'James', 'Marketing', 100000)`
- Suppose
  - T1 starts, finds tuples salary > 90000 using index and locks them
  - And then T2 executes.
  - Do T1 and T2 conflict? Does tuple level locking detect the conflict?
  - Instance of the **phantom phenomenon**
- Also consider T3 below, with Wu's salary = 90000  
`update instructor  
set salary = salary * 1.1  
where name = 'Wu'`
- Key idea: Detect “**predicate**” conflicts, and use some form of “**predicate locking**”



# Trading Consistency for Availability

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# What is Consistency?

## ■ Consistency in Databases (ACID):

- Database has a set of integrity constraints
- A consistent database state is one where all integrity constraints are satisfied
- Each transaction run individually on a consistent database state must leave the database in a consistent state

## ■ Consistency in distributed systems with replication

- Strong consistency: a schedule with read and write operations on a replicated object should give results and final state equivalent to some schedule on a single copy of the object, with order of operations from a single site preserved
- Weak consistency (several forms)



# Availability

- Traditionally, availability of centralized server
- For distributed systems, availability of system to process requests
  - For large system, at almost any point in time there's a good chance that
    - ▶ a node is down or even
    - ▶ Network partitioning
- Distributed consensus algorithms will block during partitions to ensure consistency
  - Many applications require continued operation even during a network partition
    - ▶ Even at cost of consistency



# Brewer's CAP Theorem

- Three properties of a system
  - Consistency (all copies have same value)
  - Availability (system can run even if parts have failed)
    - Via replication
  - Partitions (network can break into two or more parts, each with active systems that can't talk to other parts)
- Brewer's CAP “Theorem”: You can have at most two of these three properties for any system
- Very large systems will partition at some point
  - ➔ Choose one of consistency or availability
    - Traditional database choose consistency
    - Most Web applications choose availability
      - Except for specific parts such as order processing



# Eventual Consistency

- When no updates occur for a long period of time, eventually all updates will propagate through the system and all the nodes will be consistent
- For a given accepted update and a given node, eventually either the update reaches the node or the node is removed from service
- Known as **BASE** (**B**asically **A**vailable, **S**oft state, **E**ventual consistency), as opposed to ACID
  - **Soft state**: copies of a data item may be inconsistent
  - **Eventually Consistent** – copies becomes consistent at some later time if there are no more updates to that data item



# Availability vs Latency

- CAP theorem only matters when there is a partition
  - Even if partitions are rare, applications may trade off consistency for latency
    - ▶ E.g. PNUTS allows inconsistent reads to reduce latency
      - Critical for many applications
    - ▶ But update protocol (via master) ensures consistency over availability
  - Thus there are two questions :
    - ▶ If there is partitioning, how does system tradeoff *availability* for *consistency*
    - ▶ else how does system trade off *latency* for *consistency*



# Cloud Databases

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Data Storage on the Cloud

- Need to store and retrieve massive amounts of data
- Traditional parallel databases not designed to scale to 1000's of nodes (and expensive)
- Initial needs did not include full database functionality
  - Store and retrieve data items by key value is minimum functionality

## ► Key-value stores

- Several implementations
  - Bigtable from Google,
  - HBase, an open source clone of Bigtable
  - Dynamo, which is a key-value storage system from Amazon
  - Cassandra, from FaceBook
  - Sherpa/PNUTS from Yahoo!