



Chapter 1: Introduction

Database System Concepts, 7th Ed.

©Silberschatz, Korth and Sudarshan
See www.db-book.com for conditions on re-use



Why this course

- Computer System: Data + Process
- Many process related courses
- Less data related courses
- Data generation, transmission, storage, security, cleaning, sharing, access,
- Different data types: unstructured, semistructured, structured
- Different volumes
- Different data access requirements: write intensive, read intensive, strict transaction, analysis, training, inference
- This course deals with part of the above mentioned



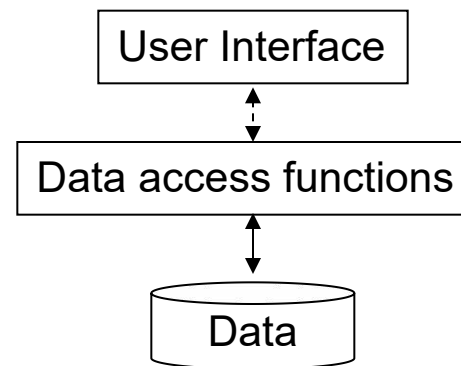
Part One: Relational Database System Concepts

- ▶ How does a relational database system work (data organization, query and update)?
- ▶ How to design a database (structure of data in the database system for a specific application)
- ▶ How to do some basic programming on the database you design



A data management scenario

- Data related problems for a course selection system
 - Where to **store** the data
 - How to **organize** the data
 - **How to access** the data
 - How to access data **efficiently**
 - How to ensure the **safety** of data
 - How to ensure the **consistency** of data
 - How to serve the increasing number of users (**concurrent** access)





Drawbacks of using file systems to store data

- Data redundancy and inconsistency
 - ▶ Multiple file formats, duplication of information in different files
- Difficulty in accessing data
 - ▶ Need to write a new program to carry out each new task
- Integrity problems
 - ▶ Integrity constraints (e.g., account balance > 0) become “buried” in program code rather than being stated explicitly
 - ▶ Hard to add new constraints or change existing ones



Drawbacks of using file systems to store data (Cont.)

- Atomicity of updates
 - ▶ Failures may leave database in an inconsistent state with partial updates carried out
 - ▶ Example: Transfer of funds from one account to another should either complete or not happen at all
- Concurrent access by multiple users
 - ▶ Concurrent access needed for performance
 - ▶ Uncontrolled concurrent accesses can lead to inconsistencies
 - Example: Two people reading a balance (say 100) and updating it by withdrawing money (say 50 each) at the same time
- Security problems
 - ▶ Hard to provide user access to some, but not all, data

Database systems offer solutions to all the above problems

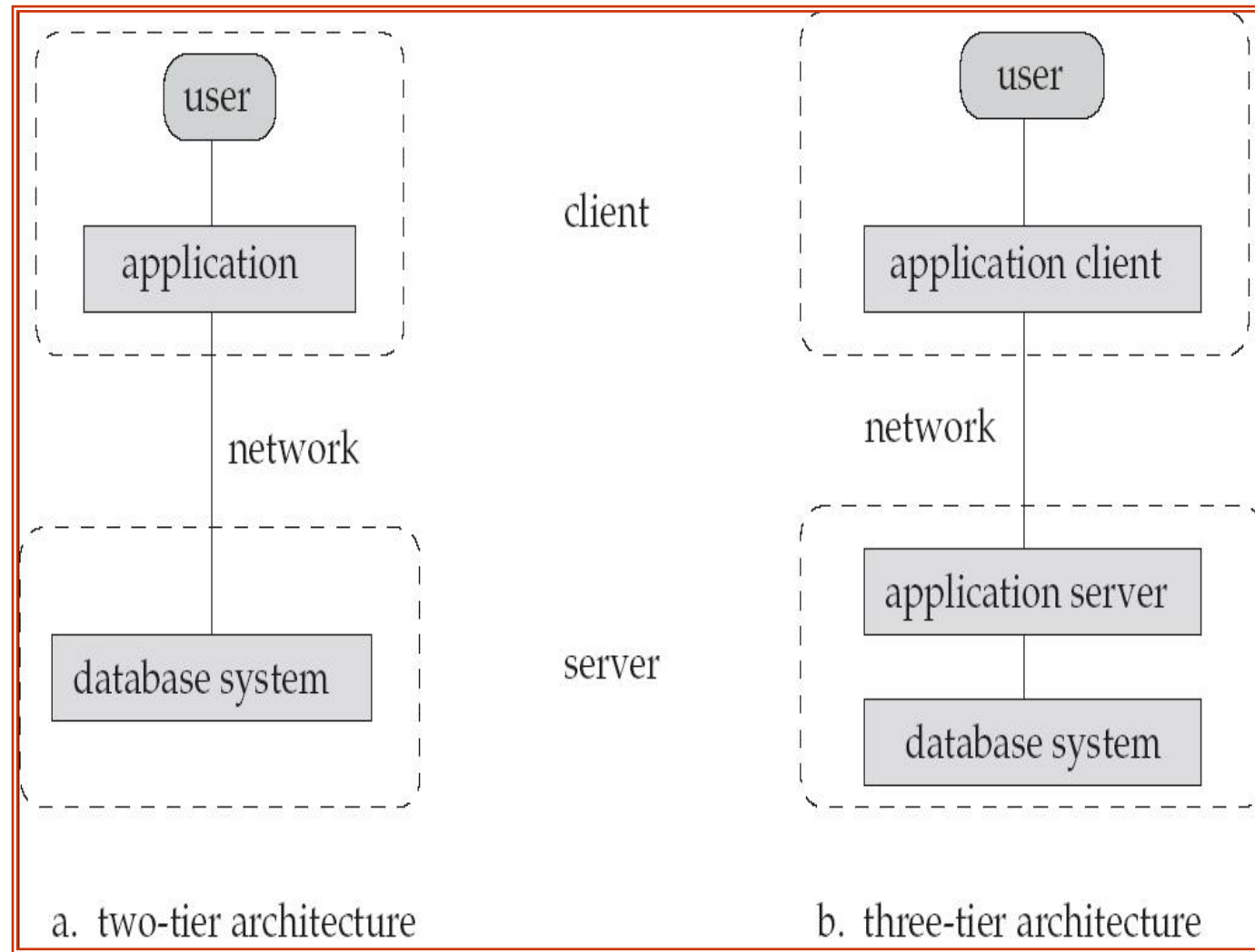


Database Management System (DBMS)

- **DBMS** contains information about a particular enterprise
 - Collection of interrelated data
 - Set of programs to access the data
 - An environment that is both *convenient* and *efficient* to use
- **Database Applications:**
 - Banking: transactions
 - Airlines: reservations, schedules
 - Universities: registration, grades, online education
 - Sales: customers, products, purchases
 - Online retailers: order tracking, customized recommendations
 - Manufacturing: production, inventory, orders, supply chain
 - Human resources: employee records, salaries, tax deductions
- **Databases** can be very large.
- Databases touch all aspects of our lives

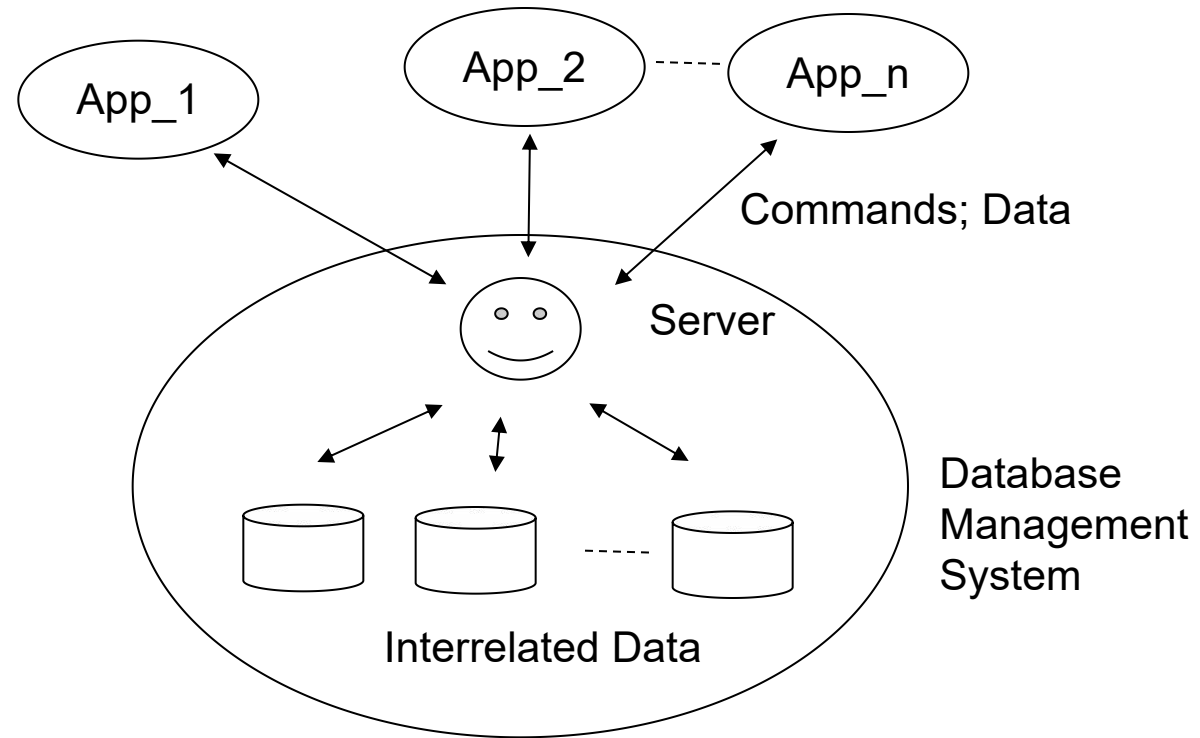


Where is database system





How does it work





Data – Different Levels

- **Physical level:** describes how a record (e.g., customer) is stored.
- **Logical level:** describes data stored in database, and the relationships among the data.

type *instructor* = **record**

```
ID : string;  
name : string;  
dept_name : string;  
salary : integer;
```

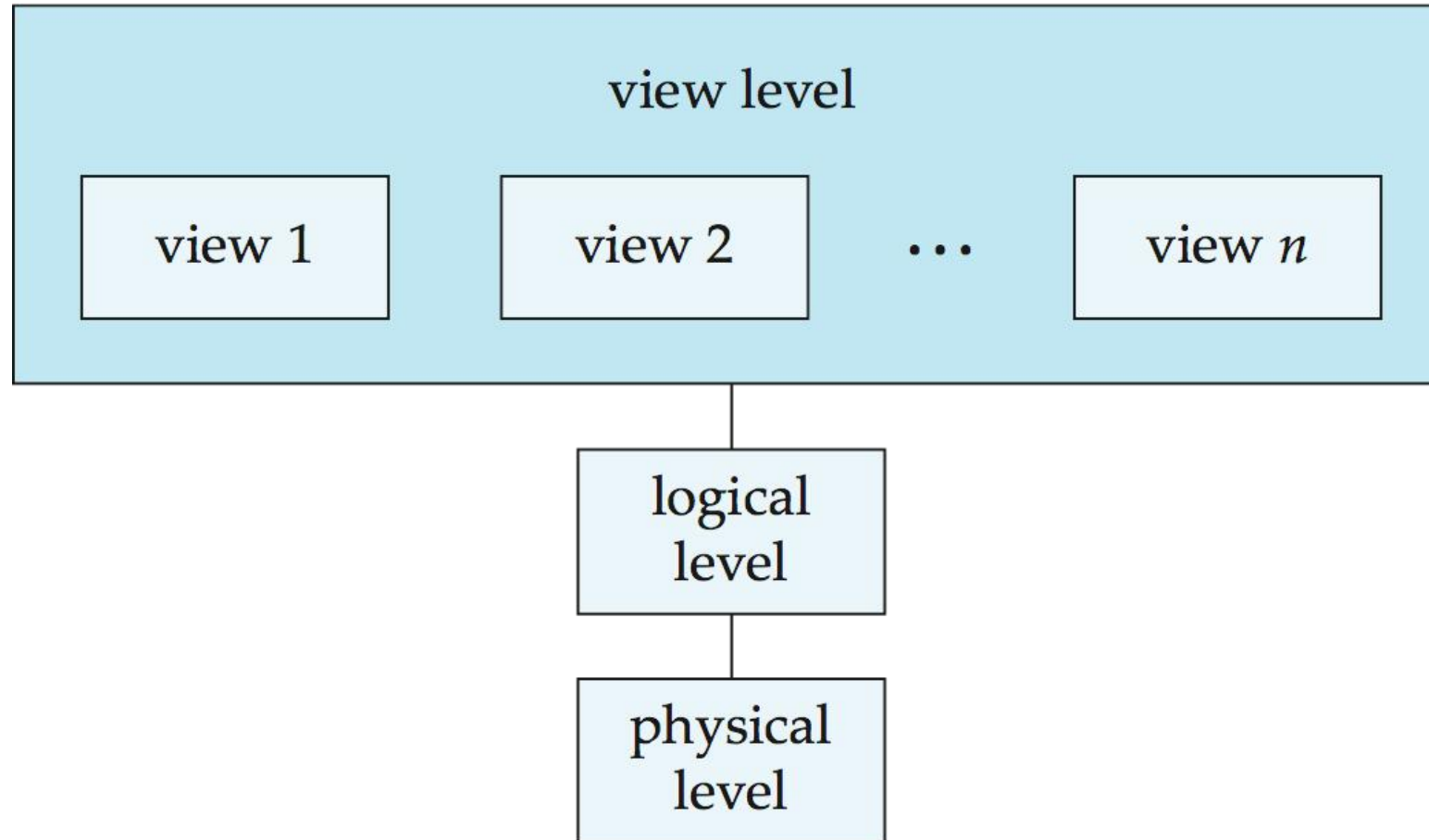
end;

- **View level:** application programs hide details of data types. Views can also hide information (such as an employee's salary) for security purposes.



Data – Different Levels (Cont.)

An architecture of data in a database system





Data Models

- A collection of tools for describing
 - Data
 - Data relationships
 - Data semantics
 - Data constraints
- Relational model, example: department(*dept_name*, *building*, *budget*)
- Entity-Relationship data model (mainly for database design)
- Object-based data models (Object-oriented and Object-relational)
- Semistructured data model (XML)
- Other older models:
 - Network model
 - Hierarchical model



Instances and Schemas

- Similar to types and variables in programming languages
- **Schema** – the logical structure of the database
 - Example: The database consists of information about a set of customers and accounts and the relationship between them
 - Analogous to type information of a variable in a program
 - **Physical schema**: database design at the physical level
 - **Logical schema**: database design at the logical level
- **Instance** – the actual content of the database at a particular point of time
 - Analogous to the value of a variable
- **Physical Data Independence** – the ability to modify the physical schema without changing the logical schema
 - Applications depend on the logical schema
 - In general, the interfaces between the various levels and components should be well defined so that changes in some parts do not seriously influence others.



Relational Model

- Relational model (Chapter 2)
- Relational algebra
- Example of tabular data in the relational model

Columns

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

Rows

(a) The *instructor* table



A Sample Relational Database & Sample Queries

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

(a) The *instructor* table

$\sigma_{dept_name="Physics"}(instructor)$

<i>dept_name</i>	<i>building</i>	<i>budget</i>
Comp. Sci.	Taylor	100000
Biology	Watson	90000
Elec. Eng.	Taylor	85000
Music	Packard	80000
Finance	Painter	120000
History	Painter	50000
Physics	Watson	70000

(b) The *department* table

$\Pi_{ID, building} (\sigma_{Instructor.dept_name="Physics"} ($
 $\sigma_{instructor.dept_name=department.dept_name}(instructor \times department)))$



Data Definition Language (DDL)

- Specification notation for defining the database schema

Example: **create table** *instructor* (
 ID **char**(5),
 name **varchar**(20),
 dept_name **varchar**(20),
 salary **numeric**(8,2))

- DDL compiler generates a set of table templates stored in a **data dictionary**
- Data dictionary contains metadata (i.e., data about data)
 - Database schema
 - Integrity constraints
 - ▶ Primary key (ID uniquely identifies instructors)
 - ▶ Referential integrity (**references** constraint in SQL)
 - e.g. *dept_name* value in any *instructor* tuple must appear in *department* relation
 - ▶ Value constraints
 - Authorization



Data Manipulation Language (DML)

- Language for accessing and manipulating the data organized by the appropriate data model
 - DML also known as query language
- Two classes of languages
 - **Procedural** – user specifies what data is required and how to get those data. Relational algebra is procedural
 - **Declarative (nonprocedural)** – user specifies what data is required without specifying how to get those data
- SQL is the most widely used query language



SQL

- **SQL**: widely used non-procedural language
 - Example: Find the name of the instructor with ID 22222

```
select  name
from    instructor
where   instructor.ID = '22222'
```
 - Example: Find the ID and building of instructors in the Physics dept.

```
select instructor.ID, department.building
from   instructor, department
where  instructor.dept_name = department.dept_name and
       department.dept_name = 'Physics'
```
- Application programs generally access databases through one of
 - Language extensions to allow embedded SQL
 - Application program interface (e.g., ODBC/JDBC) which allow SQL queries to be sent to a database
- Chapters 3, 4 and 5



Database Design

The process of designing the general structure of the database:

- Logical Design – Deciding on the database schema. Database design requires that we find a “good” collection of relation schemas.
 - Business decision – What relation schemas and attributes should we record in the database?
 - Computer Science decision – How to optimize the schema and how should the attributes be distributed among the various relation schemas?
- Physical Design – Deciding on the physical layout of the database



Database Design?

- Is there any problem with this design?

<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



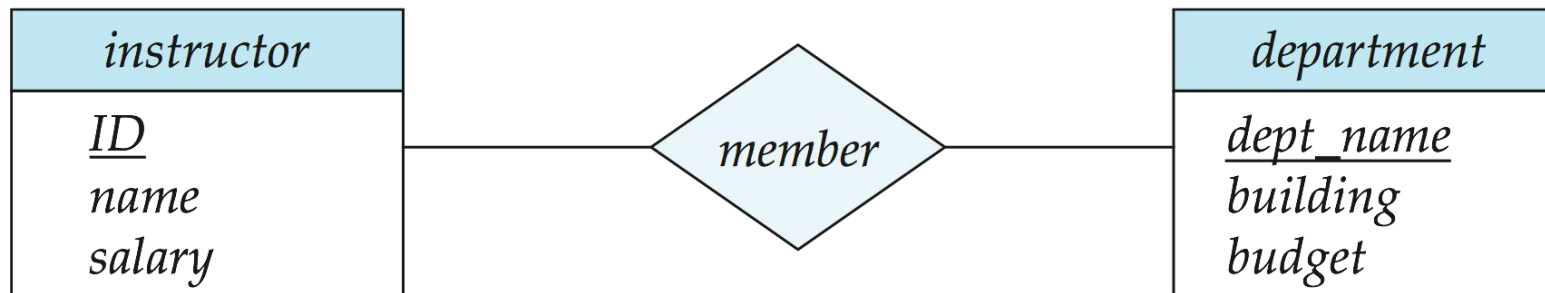
Design Approaches

- Entity Relationship Model (Chapter 6)
 - Models an enterprise as a collection of *entities* and *relationships*
 - ▶ Entity: a “thing” or “object” in the enterprise that is distinguishable from other objects
 - Described by a set of *attributes*
 - ▶ Relationship: an association among several entities
 - Represented diagrammatically by an *entity-relationship diagram*:
- Normalization Theory (Chapter 7)
 - Formalize what designs are bad, and test for them
 - Basic theory: functional dependency. Example: `student_id -> student_name`



The Entity-Relationship Model

- Models an enterprise as a collection of *entities* and *relationships*
 - Entity: a “thing” or “object” in the enterprise that is distinguishable from other objects
 - ▶ Described by a set of *attributes*
 - Relationship: an association among several entities
- Represented diagrammatically by an *entity-relationship diagram*:



Does a instructor have a dept_name?

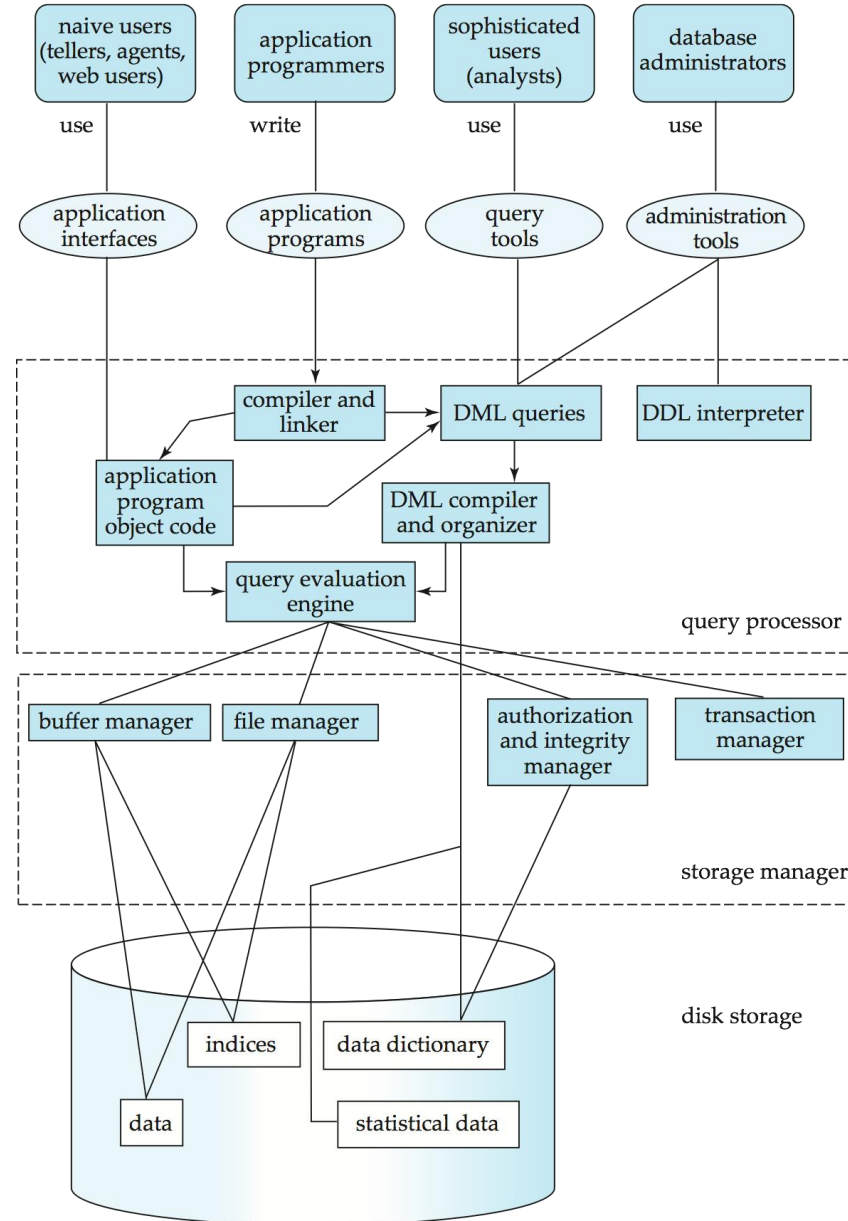


Part Two: Relational Database System Design

- How the data is physically stored and managed?
- How a query is processed?
- How a transaction is processed?
- You are required to design the storage management and query process modules plus the simple user interface



Database System Internals





Storage Management

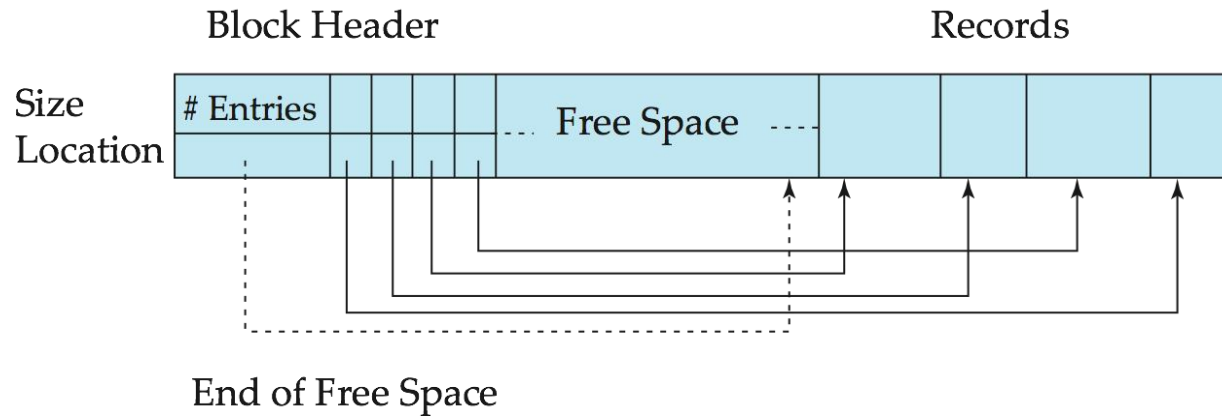
(Chapter 12,13,14)

- **Storage manager** is a program module that provides the interface between the low-level data stored in the database and the query processor of DBMS.
- The storage manager is responsible to the following tasks:
 - Interaction with the file manager
 - Efficient storing, retrieving and updating of data
- Issues:
 - Storage access (file, buffer); File organization; Indexing and hashing
 - Transaction;
 - Failure recovery
- The storage manager components include:
 - Authorization and integrity manager
 - Transaction manager
 - File manager
 - Buffer manager

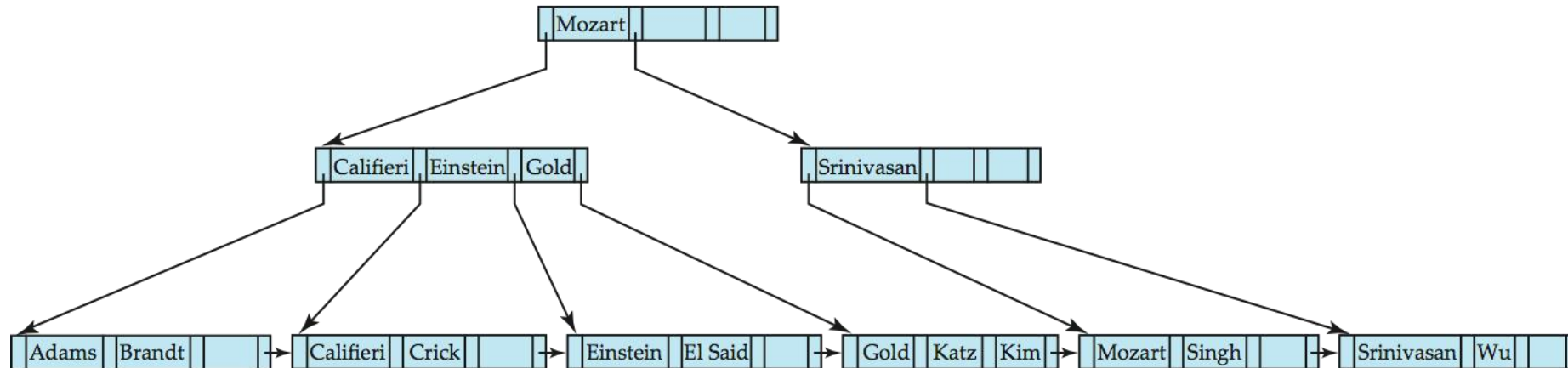


Storage Management (Cont.)

■ File organization example:



■ Index example:





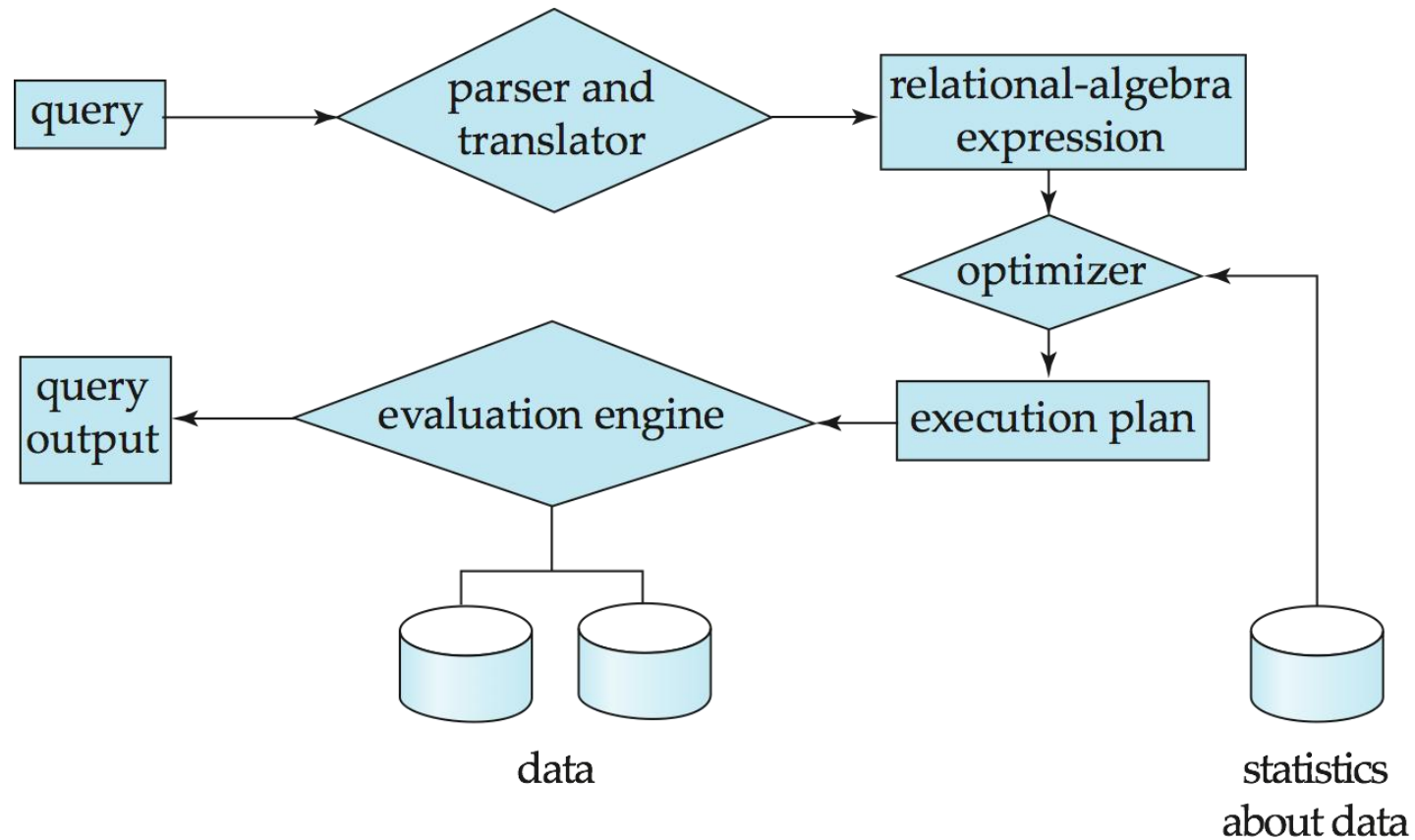
Query Processing (Chapter 15,16)

- The query processor components include:
 - DDL interpreter -- interprets DDL statements and records the definitions in the data dictionary.
 - DML compiler -- translates DML statements in a query language into an evaluation plan consisting of low-level instructions that the query evaluation engine understands.
 - ▶ The DML compiler performs query optimization; that is, it picks the lowest cost evaluation plan from among the various alternatives.
 - Query evaluation engine -- executes low-level instructions generated by the DML compiler.



Query Processing(Cont.)

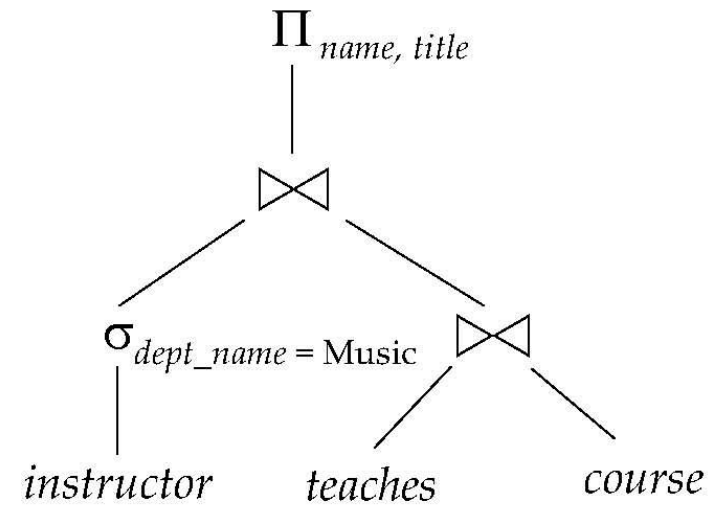
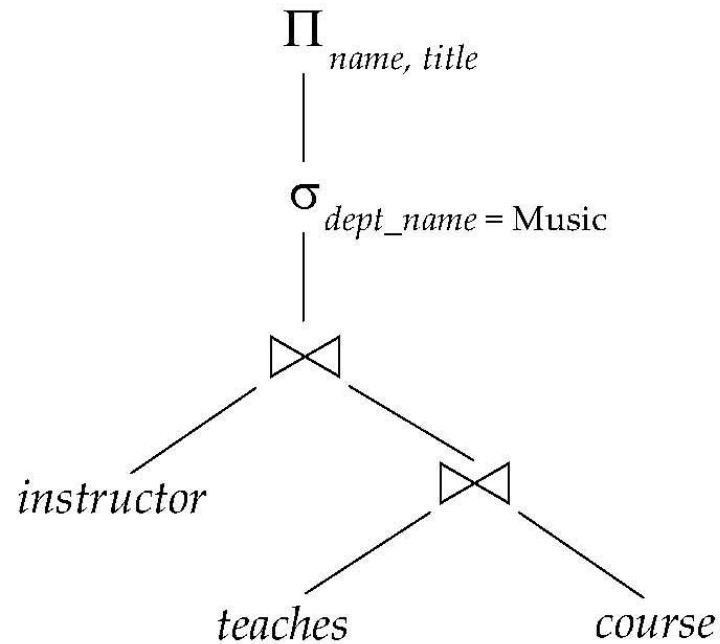
1. Parsing and translation
2. Optimization
3. Evaluation





Query Processing (Cont.)

- Alternative ways of evaluating a given query
 - Equivalent expressions
 - Different algorithms for each operation





Query Processing (Cont.)

- Different algorithms for each operation. Cost difference between a good and a bad way of evaluating a query can be enormous
 - Sequence search: $\text{Num of Blocks} * t_T + t_S$
 - Index based search: $(h_i + 1) * (t_T + t_S)$
- Need to estimate the cost of operations
 - Depends critically on statistical information about relations which the database must maintain
 - Need to estimate size of intermediate results to compute cost of complex expressions. Example, estimate number of records to find course information from takes for a specific student: $n_{\text{takes}} / V(\text{student_id}, \text{takes})$, where *takes* is the relation that stores information of students and the courses they take.



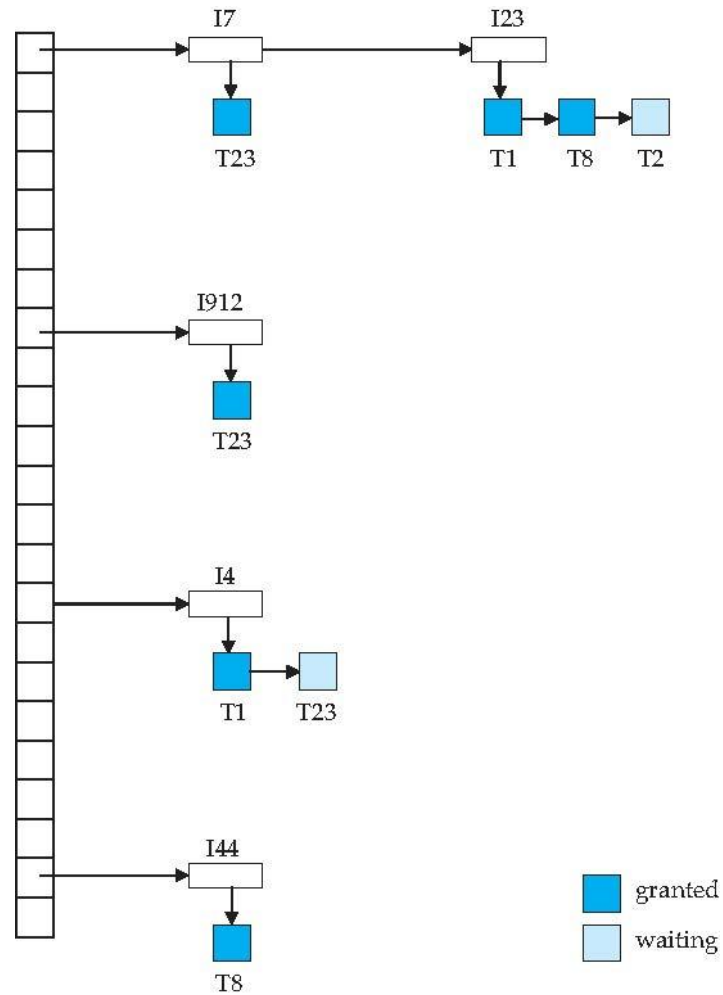
Transaction Management

Chapter 17,18,19

- What is a transaction: A **transaction** is a collection of operations that performs a single logical function in a database application
- What if the system fails?
- What if more than one user is concurrently updating the same data?
- **Transaction-management component** ensures that the database remains in a consistent (correct) state despite system failures (e.g., power failures and operating system crashes) and transaction failures.
- **Concurrency-control manager** controls the interaction among the concurrent transactions, to ensure the consistency of the database.

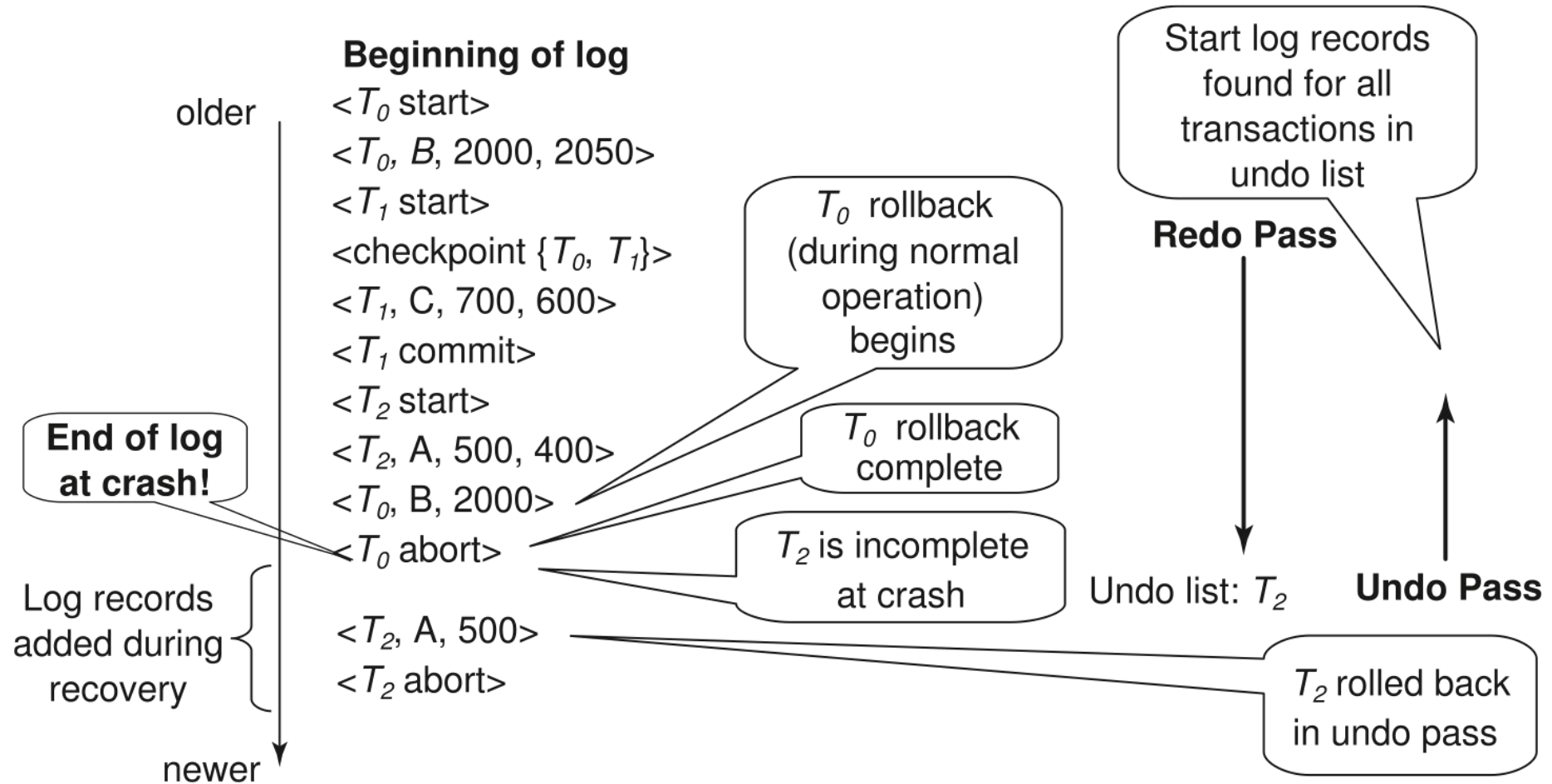


Transaction Management: Lock Based Concurrent Control



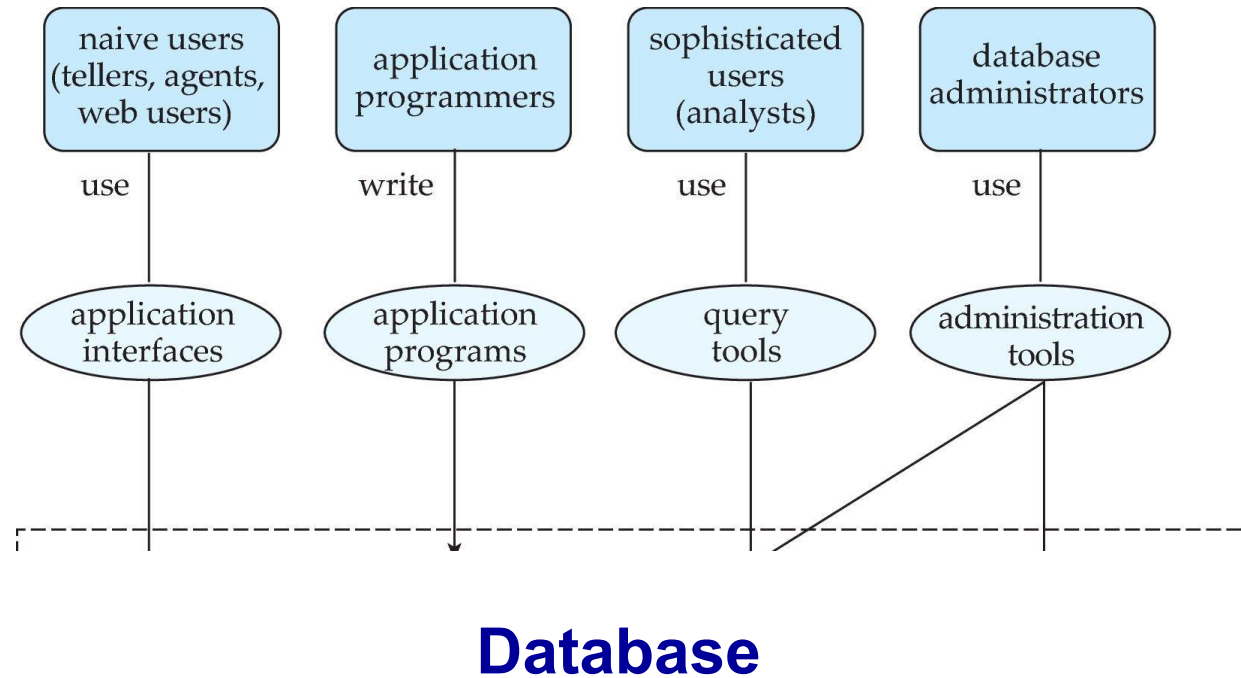


Transaction Management: Log Based Recovery





Database Users and Administrators





History of Database Systems

- 1950s and early 1960s:
 - Data processing using magnetic tapes for storage
 - ▶ Tapes provided only sequential access
 - Punched cards for input
- Late 1960s and 1970s:
 - Hard disks allowed direct access to data
 - Network and hierarchical data models in widespread use
 - Ted Codd defines the relational data model
 - ▶ Would win the ACM Turing Award for this work
 - ▶ IBM Research begins System R prototype
 - ▶ UC Berkeley begins Ingres prototype
 - ▶ Oracle releases first commercial relational database
 - High-performance (for the era) transaction processing (**OLTP**)



History (cont.)

- 1980s:
 - Research relational prototypes evolve into commercial systems
 - ▶ SQL becomes industrial standard
 - Parallel and distributed database systems
 - Object-oriented database systems
- 1990s:
 - Large decision support and data-mining applications (OLAP)
 - Large multi-terabyte data warehouses
 - Emergence of Web commerce
- 2000s:
 - NoSQL: MongoDB, Redis, HBase
 - Cloud DB: Google Spanner, Amazon Aurora, Aliyun PolarDB
- Now: AI based database or AI supporting database; Blockchain database
- It's not later products replace the earlier ones. Different products can satisfy different application requirements



End of Chapter 1

Exercise 7, 8, 9, 15



Figure 1.02

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

(a) The *instructor* table

<i>dept_name</i>	<i>building</i>	<i>budget</i>
Comp. Sci.	Taylor	100000
Biology	Watson	90000
Elec. Eng.	Taylor	85000
Music	Packard	80000
Finance	Painter	120000
History	Painter	50000
Physics	Watson	70000

(b) The *department* table

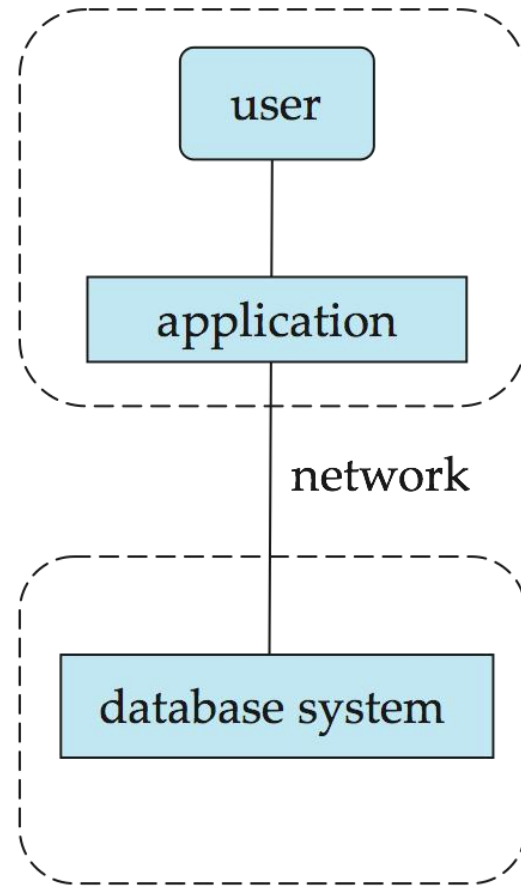


Figure 1.04

<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

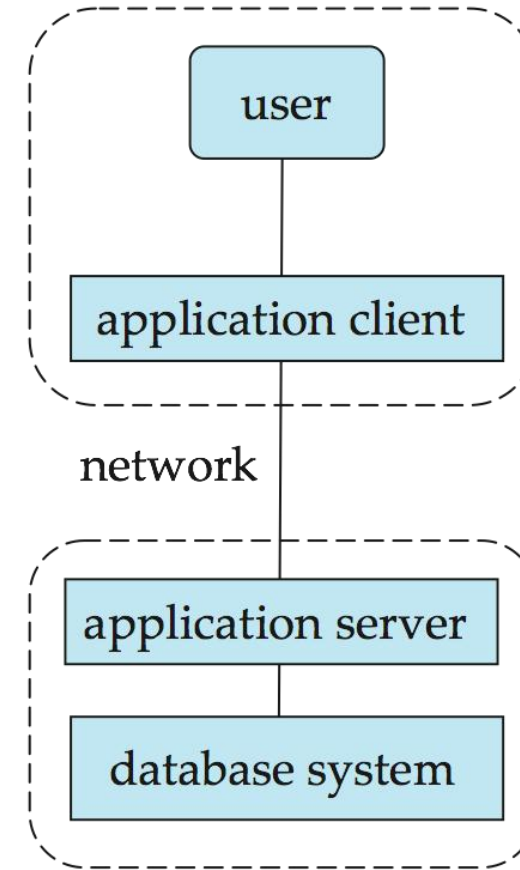


Figure 1.06



(a) Two-tier architecture

client



(b) Three-tier architecture

server