# Chapter 1: Introduction

# 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
    - 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

# University Database Example

n   Application program examples

  l   Add new students, instructors, and courses

  l   Register students for courses, and generate class rosters

  l   Assign grades to students, compute grade point averages (GPA) and generate transcripts

n   In the early days, database applications were built directly on top of file systems

# 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
- Data isolation — multiple files and formats
- 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**

# Levels of Abstraction

n   The major purpose of a database system is to provide users with an **abstract view** of the system. The system hides certain details of how data is stored and created and maintained Complexity should be hidden from database users.

n   **Physical level:** describes how a record (e.g., customer) is stored.

n   **Logical level:** describes data stored in database, and the relationships among the data.

> **type** *instructor* = **record**
>
>   *ID* : string;
>   *name* : string;
>   *dept_name* : string;
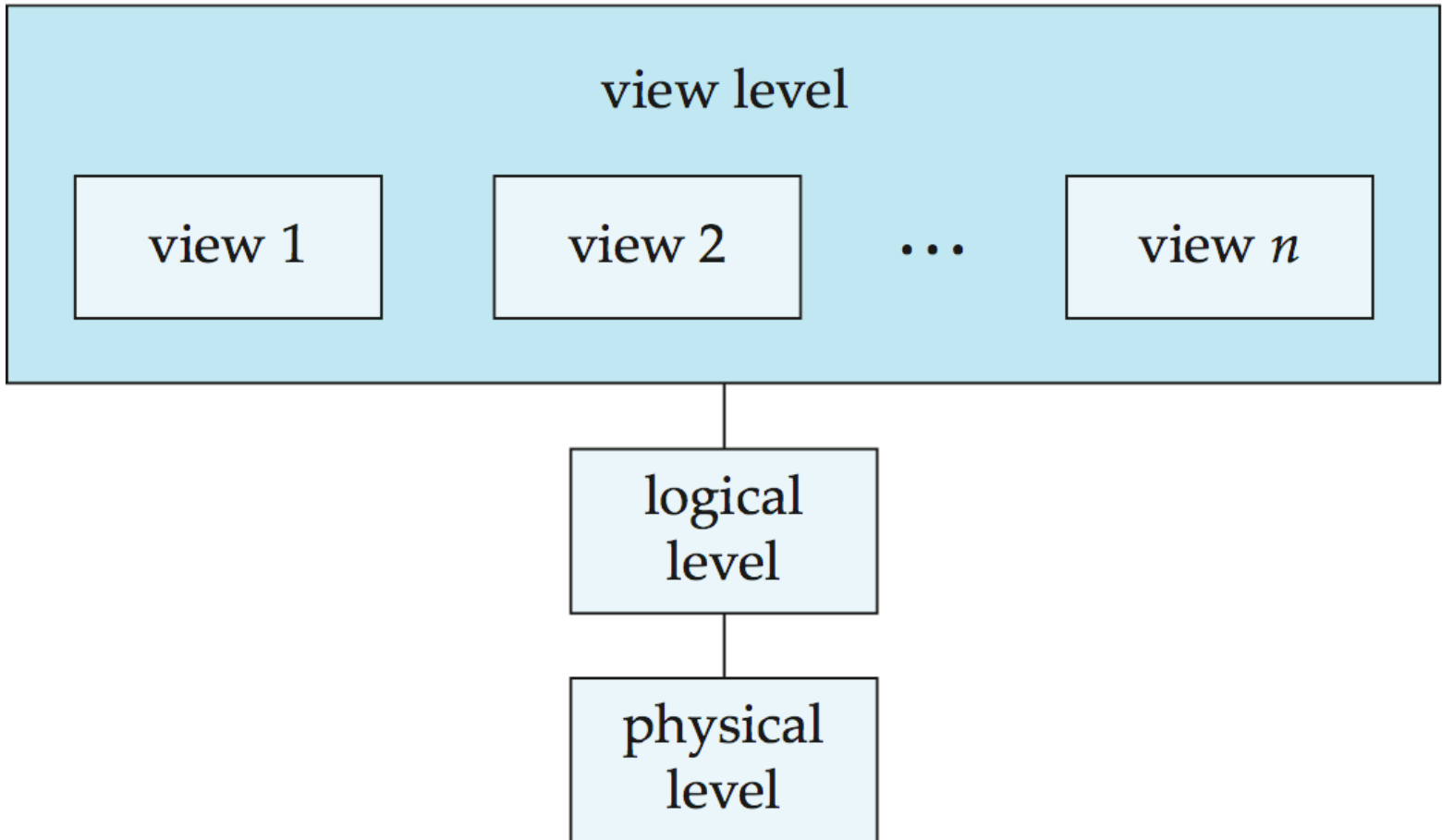>   *salary* : integer;
>
>       **end**;

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

# View of Data

An architecture for a database system

# Instances and Schemas

- **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
  - The overall design of the database is called Schema
  - 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 collection of information stored in a database at a particular moment is called instance. the actual content of the database at a particular point in 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.

# Data Models

- **Data model provides a way to describe the design of a database at the physical, logical and view level**

- A collection of conceptual tools for describing
    - Data
    - Data relationships

- Relational model: The relational model uses collection of tables to represent both data and the relationship among those data. Each table has multiple columns and each columns has a unique name

- Entity-Relationship data model (mainly for database design): This model is based on a perception of a real world that consists of a collection of basic objects called entities and relationships among these objects.

- Object-based data models (Object-oriented and Object-relational):It has seen increasing attention.

- Semistructured data model  (XML):It permits the specification of data where individual data items of the same type may have different set of attributes.

- Other older models:
    - Network model
    - Hierarchical model

# Relational Model

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

Columns

| ID | name | dept_name | salary |
|----|------|-----------|--------|
| 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

| ID | name | dept_name | salary |
|----|------|-----------|--------|
| 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



Sample E-R Diagram

| dept_name | building | budget |
|-----------|----------|--------|
| 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

# Database Language

**Data Definition Language** (**DDL**) to specify the database schema.

**Data Manipulation Language** (**DML**) to express database queries and updates.

In practice, the data-definition and data-manipulation languages are not two separate languages; instead they simply form parts of a single database language, such as the SQL language.

# Data Definition Language (DDL)

We specify a database schema by a set of definitions expressed by a special language called a data-definition language (DDL).

We specify the storage structure and access methods used by the database system by a set of statements in a special type of DDL called a **data storage and definition** language.

The DDL provides facilities to specify such constraints
**Domain Constraints**. A domain of possible values must be associated with every attribute (for example, integer types, character types, date/time types).

**Referential Integrity**. It appears in one relation for a given set of attributes also appears in a certain set of attributes in another relation (referential integrity). For example, the *dept name* value in a *course* record must appear in the *dept name* attribute of some record of the *department* relation.

# Data Definition Language (DDL)

**Authorization:**

**Read authorization**, which allows reading, but not modification, of data;

**insert authorization**, which allows insertion of new data, but not modification of existing data;

**update authorization**, which allows modification, but not deletion, of data;

**delete authorization**, which allows deletion of data.

# 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

# Data Manipulation Language (DML)

☐ A **data-manipulation language** (**DML**) is a language that enables users to access or manipulate data as organized by the appropriate data model. The types of access are:
  • Retrieval of information stored in the database.
  • Insertion of new information into the database.
  • Deletion of information from the database.
  • Modification of information stored in the database.

☐ Two classes of languages

  ☐ **Procedural** – user specifies what data is required and how to get those data

  ☐ **Declarative (nonprocedural)** – user specifies what data is required without specifying how to get those data

☐ SQL is the most widely used query language

  **A query** is a statement requesting the retrieval of information. The portion of the DML that involves information retrieval is **called 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**.**

# Database Design

**Database systems** are designed to manage large bodies of information.

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 attributes should we record in the database?

   ☐ Computer Science decision – What relation schemas should we have 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?

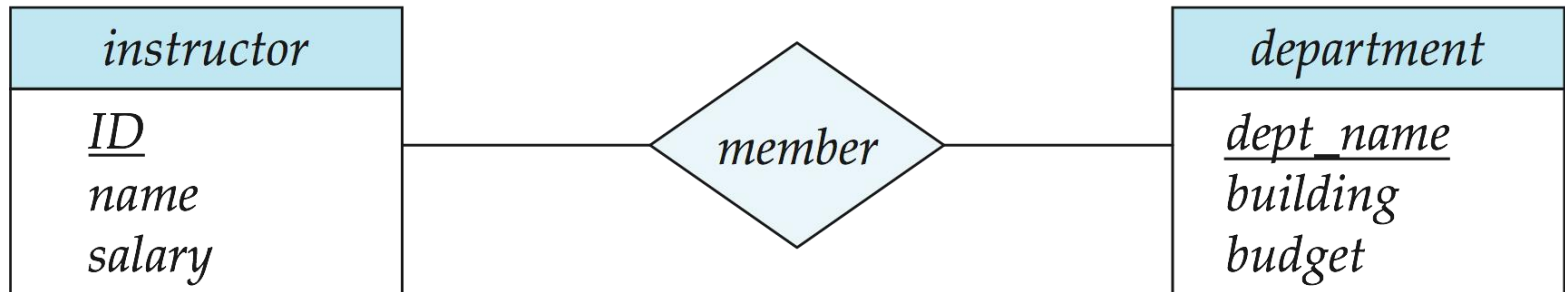| ID | name | salary | dept_name | building | budget |
|---|---|---|---|---|---|
| 22222 | Einstein | 95000 | Physics | Watson | 70000 |
| 12121 | Wu | 90000 | Finance | Painter | 120000 |
| 32343 | El Said | 60000 | History | Painter | 50000 |
| 45565 | Katz | 75000 | Comp. Sci. | Taylor | 100000 |
| 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

- Normalization Theory (Chapter 8)
  - Formalize what designs are bad, and test for them
- Entity Relationship Model (Chapter 7)
  - 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:*

# 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:*

| instructor |
| --- |
| <u>ID</u> |
| name |
| salary |

member

| department |
| --- |
| <u>dept_name</u> |
| building |
| budget |

**What happened to dept_name of instructor and student?**

# Database Engine

- A database system is partitioned into modules that deal with each of the responsibilities of the overall system.

- The functional components of a database system can be broadly divided into the

  storage manager,

  the **query processor** components,

  and the transaction management component

# Storage Management

☐ **Storage manager** is a program module that provides the interface between the low-level data stored in the database and the application programs and queries submitted to the system.

☐ The storage manager is responsible to the following tasks:

☐ Interaction with the file manager

☐ Efficient storing, retrieving and updating of data

The storage manager components include:

• **Authorization and integrity manager**, which tests for the satisfaction of integrity constraints and checks the authority of users to access data.

• **Transaction manager**, which ensures that the database remains in a consistent (correct) state despite system failures, and that concurrent transaction executions proceed without conflicts.

• **File manager**, which manages the allocation of space on disk storage and the data structures used to represent information stored on disk.

# Storage Management

☐ **Buffer manager**, which is responsible for fetching data from disk storage into main memory, and deciding what data to cache in main memory. The buffer manager is a critical part of the database system, since it enables the database to handle data sizes that are much larger than the size of main memory.

The storage manager implements several data structures as part of the physical system implementation:
• **Data files**, which store the database itself.
• **Data dictionary**, which stores metadata about the structure of the database, in particular the schema of the database.

# Storage Management

☐ **Indices**, which can provide fast access to data items. Like the index in this textbook, a database index provides pointers to those data items that hold a particular value.

For example, we could use an index to find the *instructor* record with a particular *ID*, or all *instructor* records with a particular *name*.

# Query Processor

The query processor is important because it helps the database system to simplify and facilitate access to data.

The query processor components include:
• **DDL interpreter**, which interprets DDL statements and records the definitions in the data dictionary.

• **DML compiler**, which translates DML statements in a query language into an evaluation plan consisting of low-level instructions that the query-evaluation engine understands.
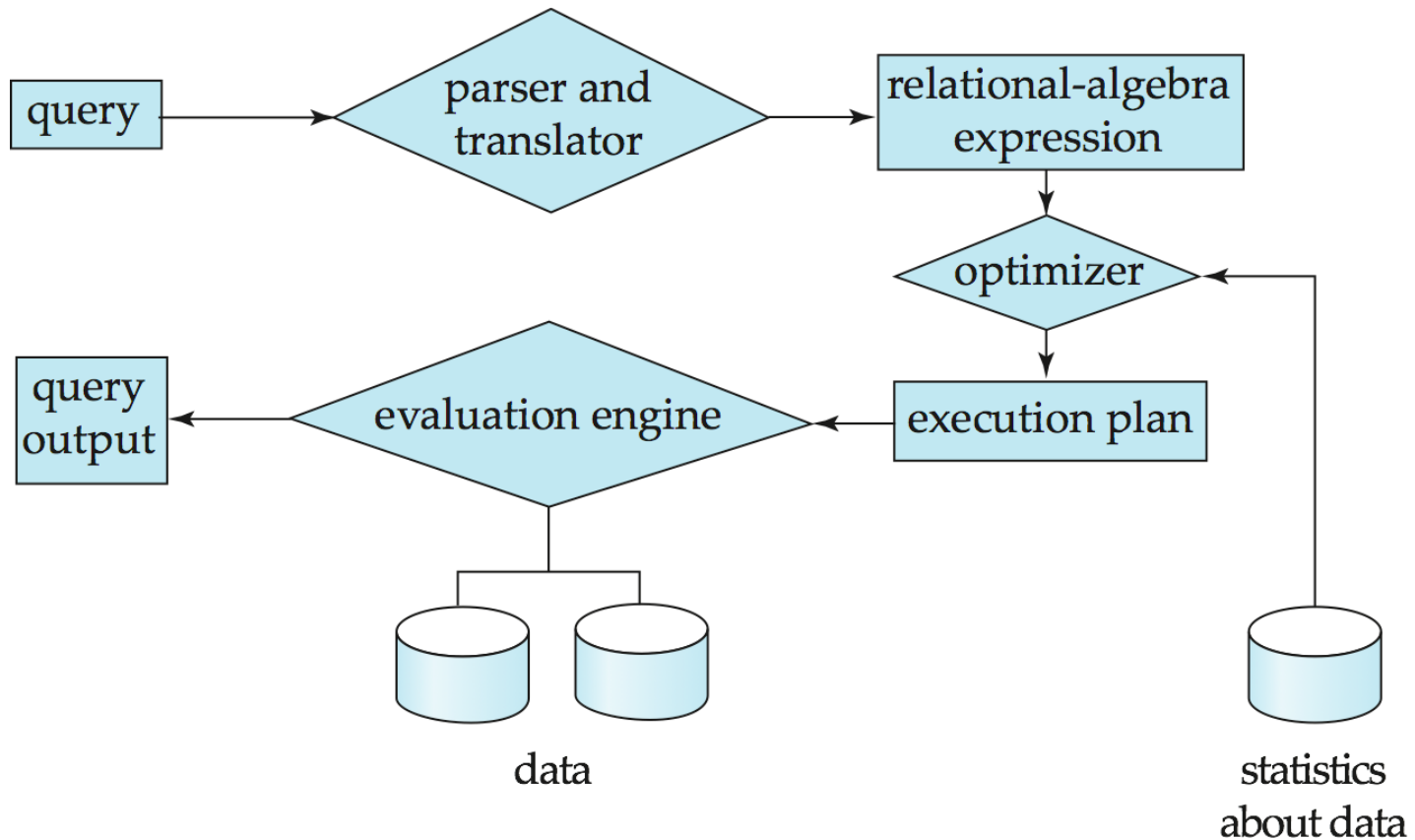
A query can usually be translated into any of a number of alternative evaluation plans that all give the same result. The DML compiler also performs **query optimization**; that is, it picks the lowest cost evaluation plan from among the alternatives.

• **Query evaluation engine**, which executes low-level instructions generated by the DML compiler.

# Query Processing

1.  Parsing and translation
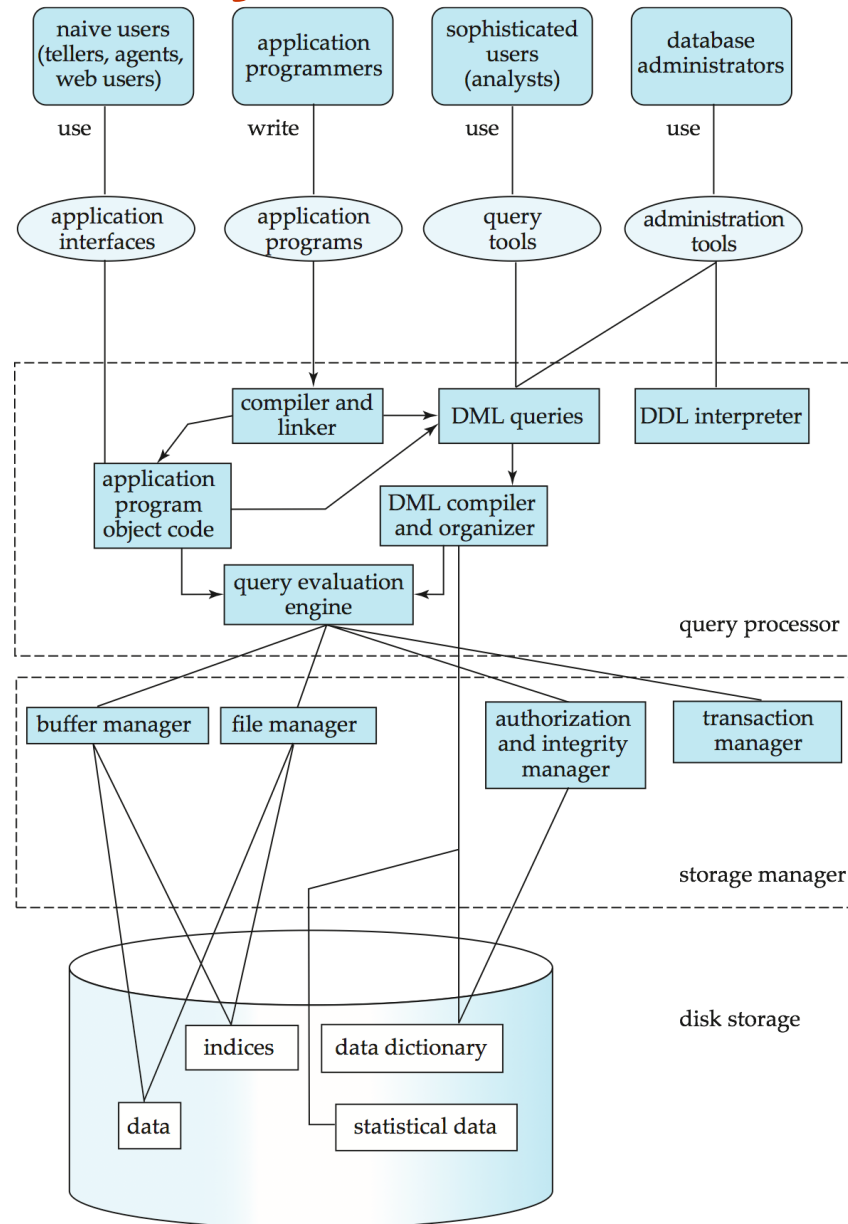2.  Optimization
3.  Evaluation

# Transaction Management

□ What if the system fails?

□ What if more than one user is concurrently updating the same data?

□ A **transaction** is a collection of operations that performs a single logical function in a database application

□ **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.

# Database System Internals



Figure shows the architecture of a database system that runs on a centralized server machine.

- The figure summarizes how different types of users interact with a database, and
- how the different components of a database engine are connected to each other.

# Database Architecture

The architecture of a database systems is greatly influenced by

 the underlying computer system on which the database is running:

- Centralized

- Client-server

- Parallel (multi-processor)

- Distributed

# 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

  ☐ High-performance (for the era) transaction processing

# 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
  - Large multi-terabyte data warehouses
  - Emergence of Web commerce
- Early 2000s:
  - XML and XQuery standards
  - Automated database administration
- Later 2000s:
  - Giant data storage systems
    - Google BigTable, Yahoo PNuts, Amazon, ..

# End of Chapter 1

# Figure 1.02

| ID | name | dept_name | salary |
|----|------|-----------|--------|
| 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

| dept_name | building | budget |
|-----------|----------|--------|
| 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

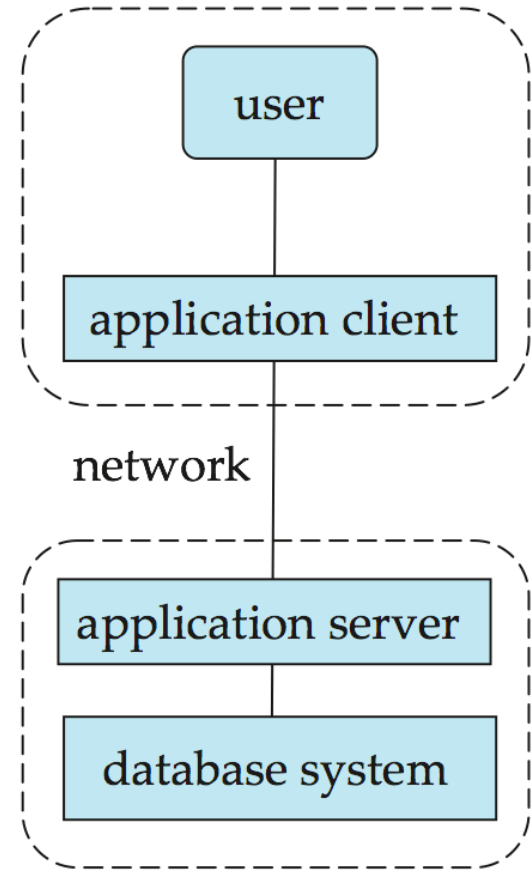| ID | name | salary | dept_name | building | budget |
|---|---|---|---|---|---|
| 22222 | Einstein | 95000 | Physics | Watson | 70000 |
| 12121 | Wu | 90000 | Finance | Painter | 120000 |
| 32343 | El Said | 60000 | History | Painter | 50000 |
| 45565 | Katz | 75000 | Comp. Sci. | Taylor | 100000 |
| 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

(b) Three-tier architecture