



# **CS425 – Fall 2019**

## **Database Organization**

### **(Graduate)**

**Yuan Hong**

**Syllabus**



# Why are Databases Important?

## □ What do Databases do?

1. Provide persistent storage
2. Efficient declarative access to data -> Querying
3. Protection from hardware/software failures
4. Safe concurrent access to data



# Who uses Databases?

- Most big software systems involve DBs!
  - Business Intelligence ⇒ e.g., IBM Cognos
  - Web based systems
  - ...
- **You!** (desktop software)
  - Your music player ⇒ e.g., Amarok
  - Your web content management system
  - Your email client
  - ...
- **Every** big company
  - Banks
  - Insurance
  - Government
  - Google, ...
  - ...





# Who Produces Databases?

- **Traditional relational database systems is big business**
  - IBM ⇒ DB2
  - Oracle ⇒ Oracle ☺
  - Microsoft ⇒ SQLServer
  - Open Source ⇒ MySQL, Postgres, ...
- **Emerging distributed systems with DB characteristics and Big Data**
  - Cloud storage and Key-value stores ⇒ Amazon S3, Google Big Table, ...
  - Big Data Analytics ⇒ Hadoop, Google MapReduce, ...
  - SQL over Distributed Platforms ⇒ Hive, Tenzing, ...





# Why are Databases Interesting (for Students)?





# Webpage and Faculty

## □ Course Info

### □ Blackboard

- ▶ Syllabus
- ▶ Announcements
- ▶ Course Materials
- ▶ Online Discussion Forum (ask questions)
- ▶ Homework & Course Project Submission

## □ Faculty

- **Yuan Hong** <http://cs.iit.edu/~yhong/>
- **Email:** [yuan.hong@iit.edu](mailto:yuan.hong@iit.edu)
- **Phone:** 312-567-5168
- **Office:** SB-216C
- **Office Hours:** Mondays, 4:30 -6:00pm (or by appointment)
- **TA(s):** TBD



# Course Objectives

- ❑ Design and model a design scenario using **relational data modeling**, which includes:
  - ❑ Analyze the design anomalies.
  - ❑ Construct **Entity Relationship Diagram**.
  - ❑ Analyze and Construct Functional Dependencies for the business rules.
  - ❑ Analyze Functional Dependencies to identify **Primary keys**.
  - ❑ Analyze and Perform **Normalization** and **Normal Forms**.
  - ❑ Define **referential integrities**.
  - ❑ Create relational database design schemas in 3-NF/BCNF for a design scenario of the size of 8-10 tables.
- ❑ Solve abstract relational language, such as **relational algebra** problems.
- ❑ Solve database transactions by using **Structured Query Language (SQL)**, used by RDBMSs.
- ❑ Explain the general concept of the **additional topics** such as: Transactions, Concurrency Control, Recovery, Structured Data and Text, and Data Warehousing.



# Tentative Schedule

Week	Dates	Topics (Chapters in the textbook)
1	8/19	Syllabus and Introduction (Chapter 1)
2	8/26	The Relational Data Model (Chapter 2), Formal Relational Query Languages (Chapter 6)
3	9/2	Labor Day (No Class) <b>Homework 1</b>
4	9/9	SQL – Introduction (Chapter 3), Homework 1 Due, Project Brainstorm
5	9/16	
6	9/23	SQL – Intermediate: Views, Integrity Constraints, Access Control (Chapter 4)
7	9/30	SQL – Advanced: APIs for SQL Access, Procedural Constructs (Chapter 5)
8	10/7	Fall Break Day (No Class) <b>Homework 2</b>
9	10/14	<b>Midterm Exam</b>
10	10/21	ER Model (Chapter 7) Homework 2 Due
11	10/28	Database Design and Normal Forms (Chapter 8)
12	11/4	Transactions (Chapter 14) <b>Homework 3</b>
13	11/11	Concurrency Control (Chapter 15) Homework 3 Due
14	11/18	Storage and Index Structures (Chapter 10 & 11)
15	11/25	Data Warehousing and Mining (Chapter 20 and extra materials)
16	12/2	Project Demo
17	TBA	<b>Final Exam</b>





# Workload and Grading

## □ Exams

- Midterm (25%) – Closed Book/Notes
- Final (30%) – Cumulative, Closed Book/Notes

## □ Homework Assignments (preparation for exams!) – 20%

- Individual assignment
- Electronic submission
- HW1 (Relational algebra)
- HW2 (SQL)
- HW3 (Database modeling)

## □ Course Project (25%)

- In groups of 2-3 students
- Given an example application (e.g., online grocery store)
  - ▶ Develop a database model
  - ▶ Derive a database schema from the model
  - ▶ Implement the application accessing the database



# Course Project

- ❑ Forming groups
  - ❑ Gather voluntarily
  - ❑ Inform me + TA as early as possible.
- ❑ Oracle Server Accounts (to be created in early September)
- ❑ Demo
  - ❑ Each group (a 10-15 minutes individual demo on 12/2)
- ❑ Submission
  - ❑ Source codes (upload a zip file on the Blackboard)
  - ❑ A short report (description of the design and implementation, some screenshots of the application, each student's contributions, etc.)
- ❑ Timeline:
  - ❑ Brainstorming on the application in September
  - ❑ Design database model (by **11/01**)
  - ❑ Derive relational model (by **11/15**)
  - ❑ Implement application and complete report (by **11/30**)



# Fraud and Late Assignments

- All work has to be original!
  - Cheating = 0 points for assignment/exam
  - Possibly E in course and further administrative sanctions
  - Every dishonesty will be reported to office of academic honesty
- Late policy:
  - -20% per day
  - No exceptions!
- Course projects:
  - Every student has to contribute in **every** phase of the project!
  - **Don't let others freeload on you hard work!**
    - ▶ Inform me or TA immediately



# Reading and Prerequisites

- **Textbook:** Silberschatz, Korth and Sudarshan
  - ***Database System Concepts, 6<sup>th</sup> edition***
  - McGraw Hill
  - New edition (7th) was released in 2019
  - ISBN: 0-07-352332-1
  - Prerequisites: CS 331 or CS401 or CS403



# For Online Students

## □ **Online Students (on campus):**

- Watch video lectures on the blackboard (will be available quickly)
- Submit homework/project on the blackboard (same as live students)
- Project:
  - ▶ will establish a platform for you to look for teammate(s)
  - ▶ coordination and source code management by your team (e.g., using GitHub)
  - ▶ one final submission would be enough
- Show up for midterm and final exams (same as live students)
- Welcome to show up in the live session



# For Online Students

## □ **Online Students (remote):**

- Watch video lectures on the blackboard (will be available quickly)
- Submit homework/project on the blackboard (same as live students)
- Project:
  - ▶ will open a platform for you to look for teammate(s)
  - ▶ coordination and source code management by your team (e.g., using GitHub)
  - ▶ one final submission would be enough
- IIT Online will contact you for midterm and final exams (assigning a nearby proctor)
- Also welcome to show up in the live session (if possible)



# Chapter 1: Introduction



# Outline

- The Need for Databases
- Data Models
- Relational Databases
- Database Design
- Storage Manager
- Query Processing
- Transaction Manager





# 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

- Application program examples
  - Add new students, instructors, and courses
  - Register students for courses, and generate class rosters
  - Assign grades to students, compute grade point averages (GPA) and generate transcripts
- 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

- **Physical level:** describes how a record (e.g., instructor) 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;
```

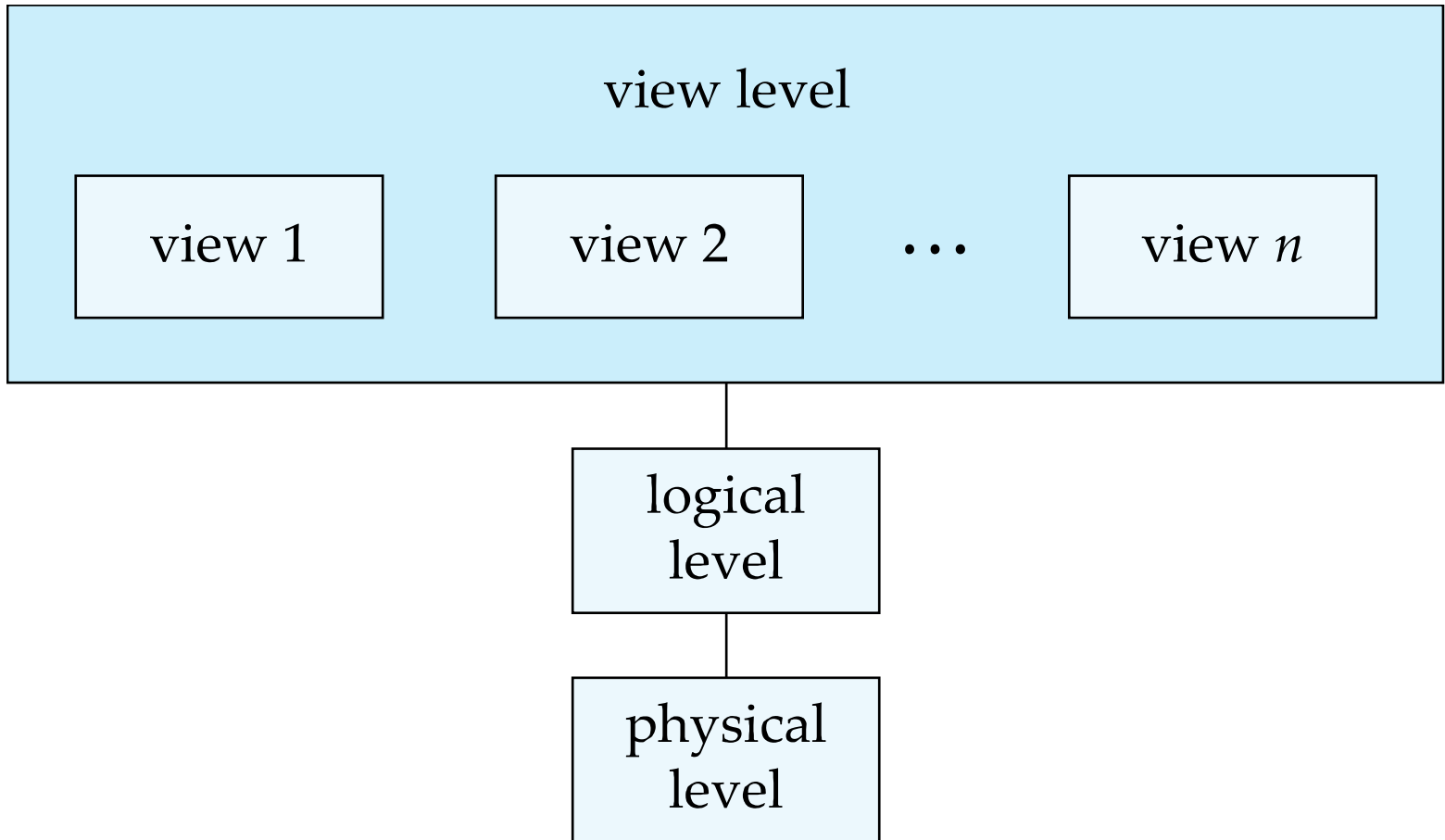
**end;**

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



# View of Data

An architecture for a database system





# Instances and Schemas

- Similar to types and variables in programming languages
- **Logical Schema** – the overall logical structure of the database
  - Example: The database consists of information about a set of customers and accounts in a bank and the relationship between them
    - ▶ Analogous to type information of a variable in a program
- **Physical schema** – the overall physical structure of the database
- **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

- ❑ A collection of tools for describing
  - ❑ Data
  - ❑ Data relationships
  - ❑ Data semantics
  - ❑ Data constraints
- ❑ Relational model
- ❑ 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





# Relational Model

- All the data is stored in various tables.
- 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

<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



# 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)
  - Authorization
    - ▶ Who can access what



# 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
  - **Pure** – used for proving properties about computational power and for optimization
    - ▶ Relational Algebra
    - ▶ Tuple relational calculus
    - ▶ Domain relational calculus
  - **Commercial** – used in commercial systems
    - ▶ SQL is the most widely used commercial language



# SQL

- ❑ The most widely used commercial language
- ❑ SQL is NOT a Turing machine equivalent language
- ❑ To be able to compute complex functions SQL is usually embedded in some higher-level language
- ❑ Application programs generally access databases through one of
  - ❑ Language extensions to allow embedded SQL
  - ❑ Application program interface (e.g., ODBC/JDBC) which allows SQL queries to be sent to a database



# 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 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 (Cont.)

- Is there any problem with this relation?

<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

- Need to come up with a methodology to ensure that each of the relations in the database is “good”
- Two ways of doing so:
  - Entity Relationship Model (Chapter 7)
    - ▶ Models an enterprise as a collection of *entities* and *relationships*
    - ▶ Represented diagrammatically by an *entity-relationship diagram*:
  - Normalization Theory (Chapter 8)
    - ▶ Formalize what designs are bad, and test for them





# Object-Relational Data Models

- ❑ Relational model: flat, “atomic” values
- ❑ Object Relational Data Models
  - ❑ Extend the relational data model by including object orientation and constructs to deal with added data types.
  - ❑ Allow attributes of tuples to have complex types, including non-atomic values such as nested relations.
  - ❑ Preserve relational foundations, in particular the declarative access to data, while extending modeling power.
  - ❑ Provide upward compatibility with existing relational languages.



# XML: Extensible Markup Language

- ❑ Defined by the WWW Consortium (W3C)
- ❑ Originally intended as a document markup language not a database language
- ❑ The ability to specify new tags, and to create nested tag structures made XML a great way to exchange **data**, not just documents
- ❑ XML has become the basis for all new generation data interchange formats.
- ❑ A wide variety of tools is available for parsing, browsing and querying XML documents/data



# Database Engine

- ❑ Storage manager
- ❑ Query processing
- ❑ Transaction manager



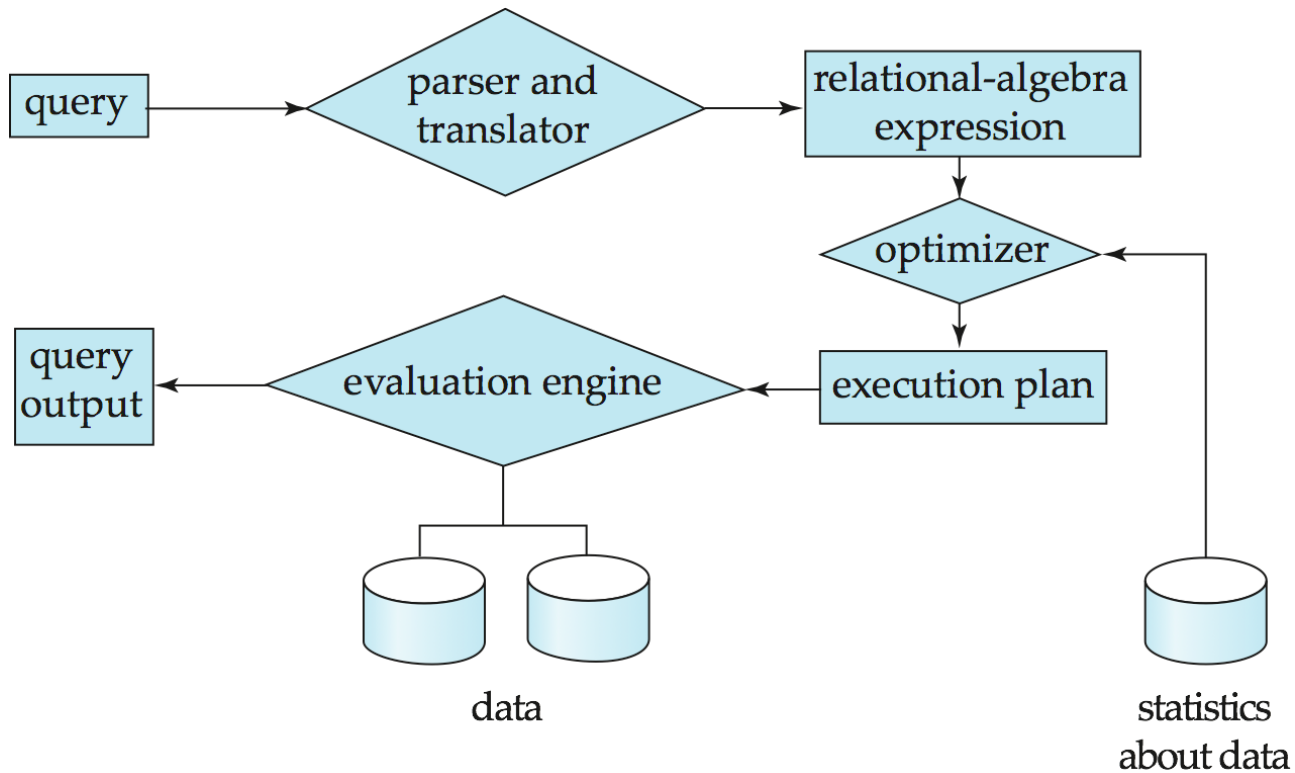
# Storage Management

- **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 OS file manager
  - Efficient storing, retrieving and updating of data
- Storage manager implements several data structures as part of the physical system implementation:
  - Data files
  - Data dictionary
  - Indices



# Query Processing

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
- Cost difference between a good and a bad way of evaluating a query can be enormous
- Need to estimate the cost of operations
  - Depends critically on **statistical information** about relations which the database must maintain
  - Need to estimate statistics for intermediate results to compute cost of complex expressions

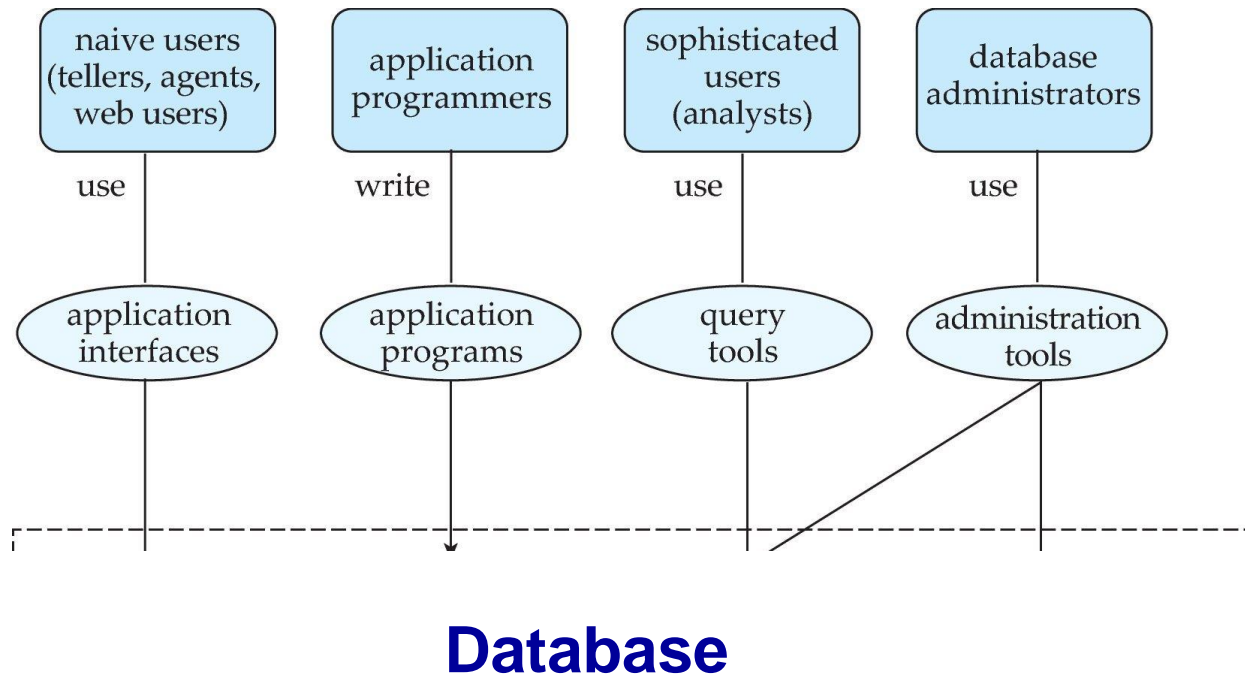


# Transaction Management

- 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.



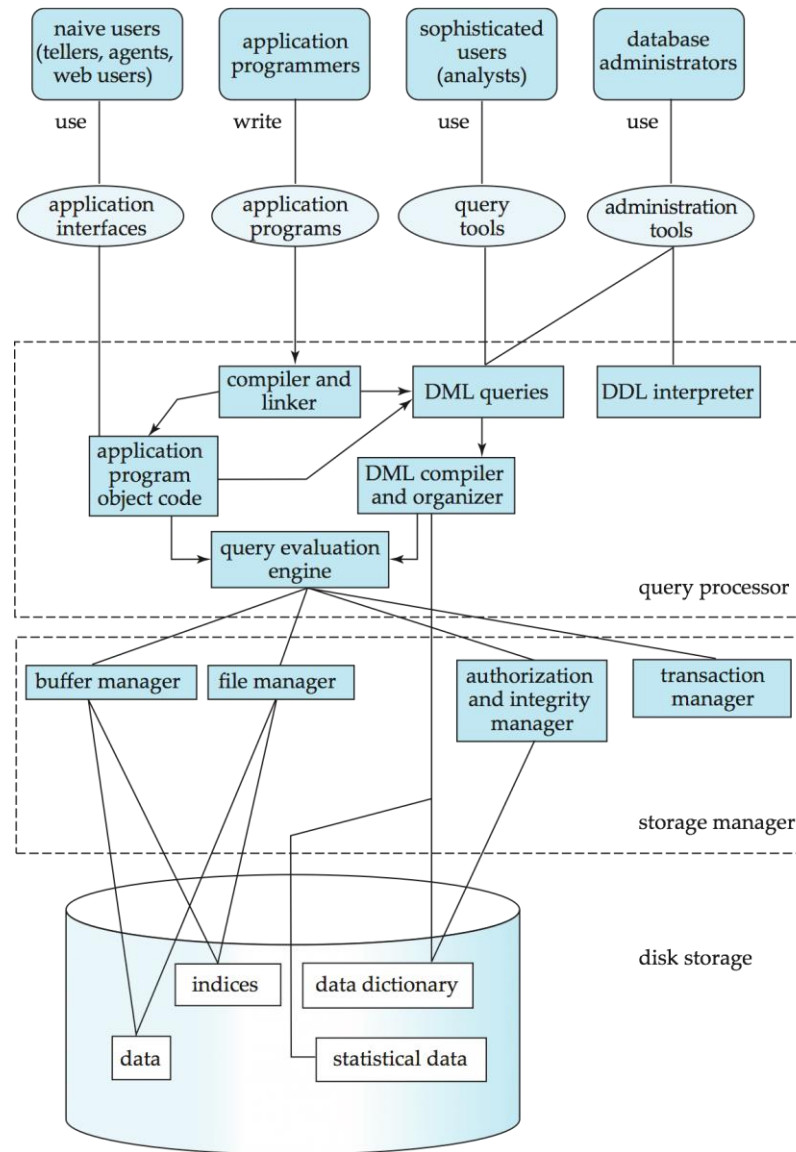
# Database Users and Administrators







# Database System Internals





# Database Architecture

The architecture of a database system 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