# EECS 484: Database Management Systems

Instructor: H. V. Jagadish

Email: jag@umich.edu

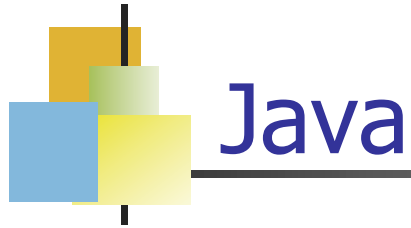Office Hours: MW 10.30-11, BBB 4601

# Course Outline – EECS 484

- GOAL: Basic introduction to database management systems.

- Two perspectives:

  - External (*Database user*)

    - Data models, ER model, relational model, SQL, database design …

    - Java/JDBC Project: Common platform for building database applications

  - Internal (*Database implementer*)

    - File organizations, access methods, sorting, concurrency control, recovery, …

    - Minirel Project: Build components of a Relational Database System

- Textbook "Database Management Systems", by Raghu Ramakrishnan & Johannes Gehrke. 3rd ed.

# Java

- Databases are most often accessed via SQL.

- But SQL is usually embedded in, and called from, a traditional programming language.

- Java is the most common choice, and so the one we have chosen for this course.

- Discussion this week is a tutorial on: Intro to Java for C++ programmers

# Groups

- Total of four group projects.
- Group of size 2
- Same partner for all projects
- Start looking for partners now!
- Register your group next week by following the link we provide.
  - If you are not in a group by next Friday, you will be assigned a partner at random.

# Project Grading

- Mostly autograder, some human.
- Limited number of submissions, even for autograded portion.
  - Make sure to test extensively.
- Same score for both partners.

# Discussion Sections

- Not optional!

- Project covered in the discussion section.

- Exams may have project-related questions.

- 5 sections on Fridays – identical content, but different instructor.

EECS 484

# Academic Honesty

CoE Honor Code for all students.

Specifics on course website.

Can discuss, but cannot copy.

Questions? – Ask me first!

# Course Policies

- Projects
  - Due by 11.55 p.m. on due date
  - 2 free late days <u>total</u> for all projects.
  - 1% course grade for each late day (or part thereof) used beyond the free day.
  - Up to 4 late days allowed per project.

- Assignments
  - Due by 11.55 p.m. on due date
  - No late submissions accepted.
  - Worst assignment dropped from total.

# Course Grading

| | |
|---|---|
| First Exam | 25% |
| Second Exam | 25% |
| 4 of 5 written homework assignments [each worth 2.5%] | 10% |
| Four projects [each worth 10%] | 40% |

**No make up exams**

EECS 484

# Karma Points

You can earn plus/minus 2% karma points.

You earn positive points by:
- Helpful posts on piazza
- Good questions in class
- Etc.

You earn negative points by:
- Posting duplicate questions on piazza without checking first
- Disrupting class
- Etc.

# Exams

- Two exams:
- One midterm (Oct 23) and one final (Dec 15).
- Non-cumulative.
- No alternate exams, no make-up exams.
- Closed book but a one page, one-side, handwritten cheat sheet is permitted.

# This week

- There is discussion.
  - Java basics for a C++ programmer.
  - (You will need this for project 2).
  - No need to go if you know Java.
- No office hours.
  - Regular office hours start next week.

# Lectures

- Notes are posted on canvas the night before.

- Sometimes updated after lecture.
  - To fix errors
  - To add clarifications

- Video recordings will be posted on canvas, usually a day or two after.

# What Is a DBMS?

- DBMS = Database Management System
- Database: Large, integrated collection of data.
- Models some real-world *enterprise*
  - Entities (e.g., students, courses)
  - Relationships (e.g., Lisa Simpson is taking EECS 484)
- DBMS: a software package designed to store and manage databases

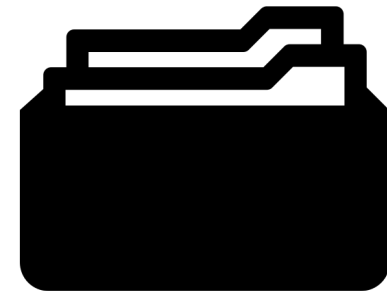# Old-time Solution: Sorted Student Folders

- Advantages?

- Disadvantages?
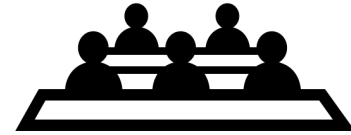
# Old-time Solution: Sorted Student Folders

- Advantages?
  - cheap

- Disadvantages?
  - Large physical footprint
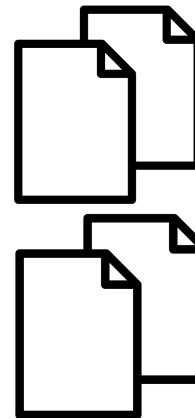  - No sharing
  - No ad-hoc queries
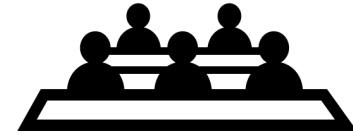
# Other Solution: Flat Files

- Access?
  - using programs in C, Java, Python etc.
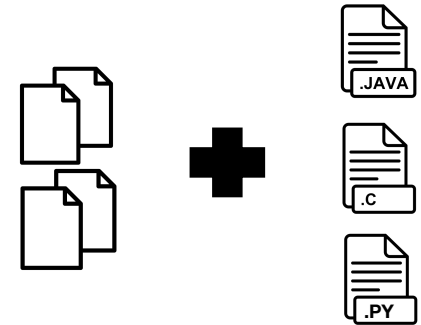- Layout for the student records?

# Other Solution: Flat Files

- ## Access?
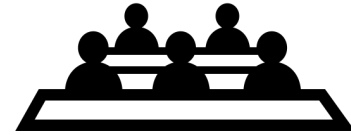  - using programs in C, Java, etc.
- ## Layout for the student records?

**CSV:**
Brown, Lisa, lbrown, db, A, os, B
Smith, Bart, bsmith
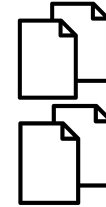Tompson, Mary, mtom, vis, B+, db, A-
…

…

# Other Solution: Flat Files

- ## Access?
  - using programs in C, Java, etc.

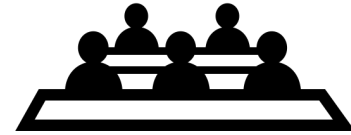- ## Layout for the student records?

**Multiple files:**

| | |
|---|---|
| Brown, Lisa, lbrown | lbrown, db, A |
| Smith, Bart, bsmith | lbrown, os, B |
| Tompson, Mary, mtom | mtom, vis, B+ |
| … | mtom, db, A- |
| … | … |
| | … |

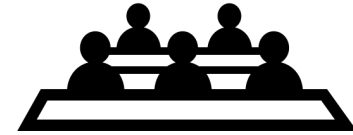# Other Solution: Flat Files
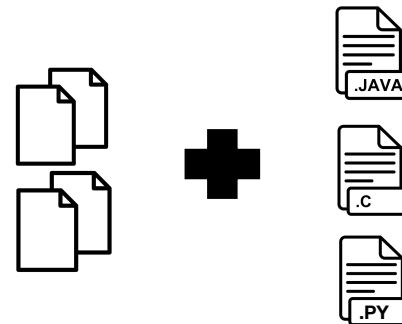
- Problems?

# Other Solution: Flat Files

- Problems?
  - Inconvenient access to data
    - requires programming experience and knowledge of file layout
  - Data redundancy
  - Integrity problems
  - Atomicity problems (concurrent access issues)
  - Security problems

# Why use a DBMS?

- It solves ALL these problems!
  - Data independence
    - Apps need a view of the data, not info about internal representation and storage
  - Efficient storage and access
  - Centralized data administration
  - Data integrity and security
  - Concurrent access, recovery from crashes
  - Reduced application dev time
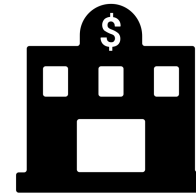
# Who uses a DBMS?

# Who uses a DBMS?

- Everyone!
  - Your bank
  - Your university
  - Your coffee shop
  - Your favorite hotel
  - Your favorite website
  - Your phone
  - Your government

- How many databases have you used so far today?

# Why Study Databases?

- Data is useless without the tools to extract information (queries)
  - "Optimal" pricing of an airline ticket

- Datasets increasing in diversity and volume
  - Websites, digital libraries, interactive video, human genome project, mobile applications

- Databases touch most of CS
  - OS, languages, theory, AI, multimedia, logic, …

# Data Models

- Data model: a collection of concepts for describing data.

- Schema: a description of a particular collection of data, using a given data model.

- Relational model: the most widely-used model today.
  - Data model: Database is a collection of relations
    A relation is a table with rows and columns.
  - Every relation has a schema, which describes the columns (also called the fields or attributes).

- Entity-Relationship (ER) model: A "semantic" data model, i.e. a higher-level more user-intuitive model
  - A (relational) DBMS only understands the relational model
    ➔ Must translate an ER schema to a relational schema

# Relational and Other Data Models

- **DBMS using the relational DM** ('70s-'80s)
  - IBM DB2
  - Informix
  - Oracle
  - Sybase
  - Microsoft Access
  - Tandem
  - Teradata
  - ...

- **Other data models**
  - Hierarchical (mid '60s-'70s)
    - IBM IMS
  - Network ('70s)
    - IDMS, IDS
  - Object-oriented (~'90s)
    - ObjectStore
  - Object-relational (relational model + object DB concepts)
    - Oracle
  - ...

# Relational (Data) Model

- The most widely-used model today

- **Data model** = a collection of concepts for describing data

  - A collection of relations

  - Relation = set of records – think of it as a table with rows and columns

Students

| sid | name | login | age |
|-----|------|-------|-----|
| 13 | Lisa | lsimp | 40 |
| 41 | Bart | bart | 20 |

Courses

| cid | cname | cred. |
|-----|-------|-------|
| E-484 | EECS484 | 4 |
| E-584 | EECS584 | 3 |

Enrolled

| sid | cid | grade |
|-----|-----|-------|
| 41 | E-484 | A- |
| 13 | E-584 | A+ |

# Relational (Data) Model

- **Schema** = a description of data in terms of a data model
  - Every relation has a schema
  - Specifies the name of the relation, the name and type of the columns (or fields or attributes)
  - Each row also called a tuple or a record

  Students(sid:`string`, name:`string`, login:`string`, age:`integer`)
  Courses(cid:`string`, cname:`string`, credits:`integer`)
  Enrolled(sid:`string`, cid:`string`, grade:`string`)

Students

| sid | name | login | age |
|-----|------|-------|-----|
| 13 | Lisa | lsimp | 40 |
| 41 | Bart | bart | 20 |

Courses

| cid | cname | cred. |
|-----|-------|-------|
| E-484 | EECS484 | 4 |
| E-584 | EECS584 | 3 |

Enrolled

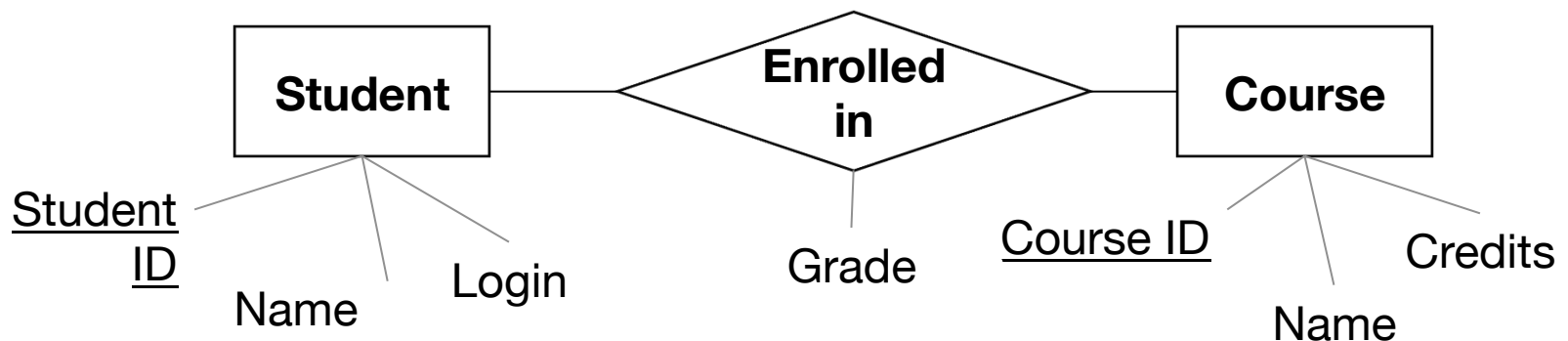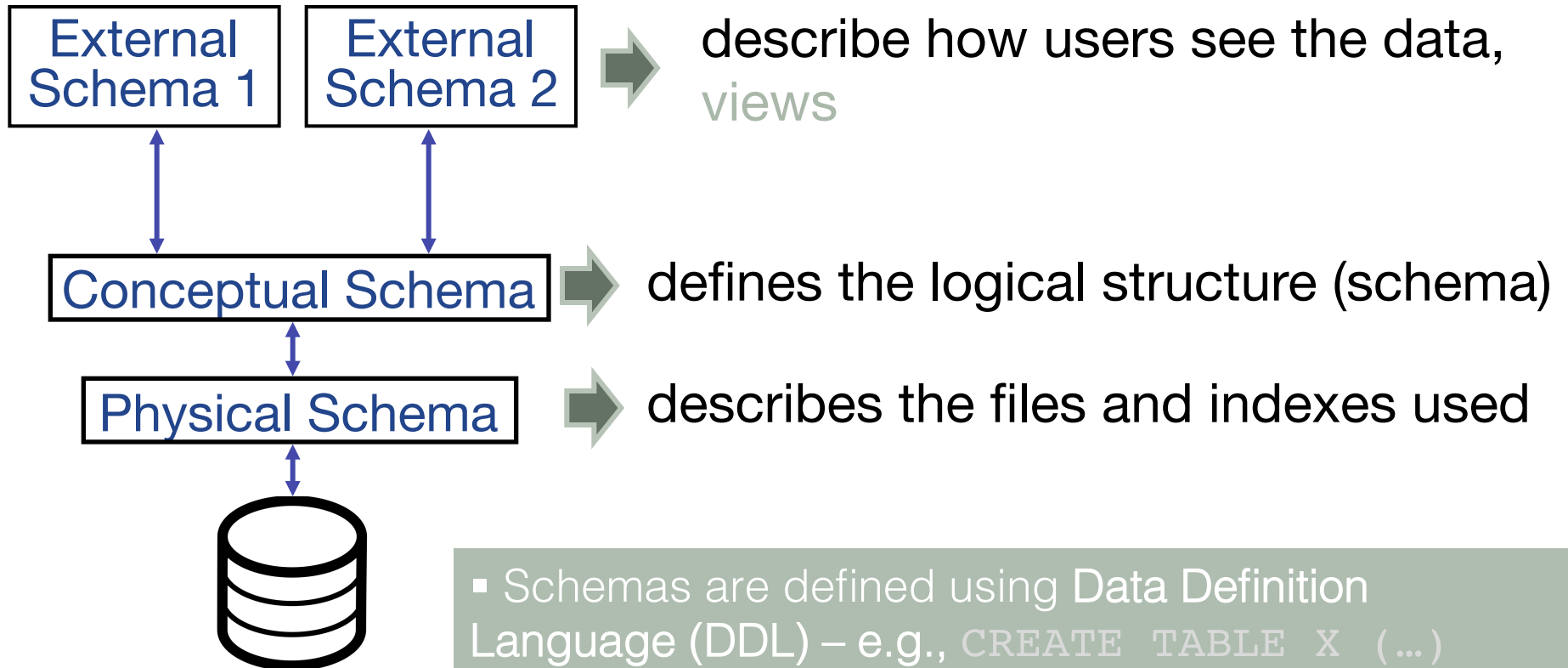| sid | cid | grade |
|-----|-----|-------|
| 41 | E-484 | A- |
| 13 | E-584 | A+ |

# Entity-Relationship (ER) Model

- A "semantic" data model
  - a higher-level, more user-intuitive model
  - A (relational) DBMS understands the relational model
    - ➔ Must translate an ER schema to a relational schema
- Entity-Relationship diagram:
  - Entities: Student, Course
  - Relationship: Enrolled_in

```
┌──────────┐      ╱◇╲              ┌──────────┐
│ Student  │────╱ Enrolled ╲──────│  Course  │
│          │    ╲   in    ╱        │          │
└──────────┘      ╲◇╱              └──────────┘
```

Student ID    Name    Login         Grade     Course ID    Name    Credits

# Levels of Abstraction

| External Schema 1 | External Schema 2 | ➧ describe how users see the data, views |

defines the logical structure (schema)

**Conceptual Schema** ➧ defines the logical structure (schema)
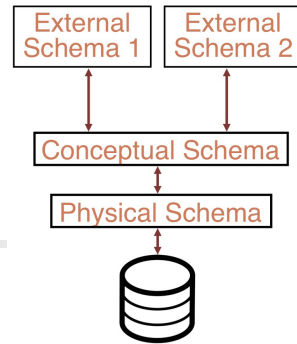
**Physical Schema** ➧ describes the files and indexes used

- Schemas are defined using **Data Definition Language (DDL)** – e.g., `CREATE TABLE X (...)`
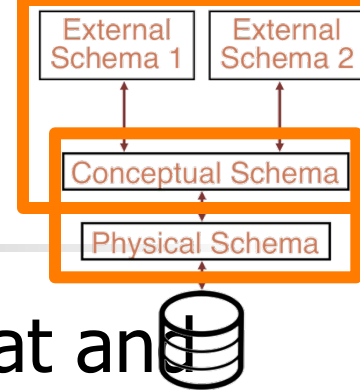- Data is modified/queried using **Data Manipulation Language (DML)** – e.g., `SELECT FROM X WHERE ...`

# Example

- Conceptual schema (1):
  - Students(sid:`string`, name:`string`, login:`string`, age:`integer`)
  - Courses(cid:`string`, cname:`string`, credits:`integer`)
  - Enrolled(sid:`string`, cid:`string`, grade:`string`)

- Physical schema (1):
  - Relations stored as unordered files.
  - Index on first column of Students.

- External Schema (≥ 1):
  - View: Course_info(cid:`string`, enrollment:`integer`)
  - View: Class_rank(sid:`string`, gpa:`real`, rank:`integer`)

# Data Independence



- Applications insulated from data format and storage details

- Logical data independence:  Protection from changes in *logical* structure of data
  - External / Conceptual schemas

- Physical data independence: Protection from changes in *physical* structure of data
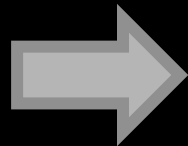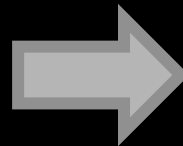  - Conceptual / Physical schemas

# CYU

- Which of these are more suitable for storing in a DBMS rather than files in an OS?

  (a) Grades for students at the university

  (b) Source code for a program
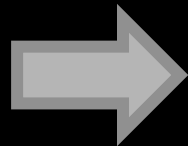
  (c) Contents of a textbook

Think → Pair → Share

# CYU

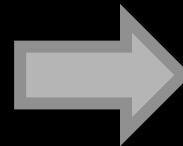- Let's say UM provides you access to a relational table that gives just your grades in various courses. Does that relation represent:

    a) An external schema?

    b) A conceptual schema?

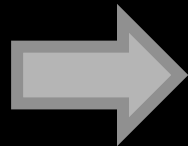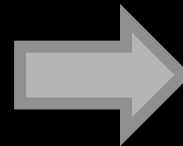    c) A physical schema?

Think → Pair → Share

# CYU

- The relational table with student grade information is very large and stored on multiple servers for performance. Does the storage scheme represent:

  a) An external schema?
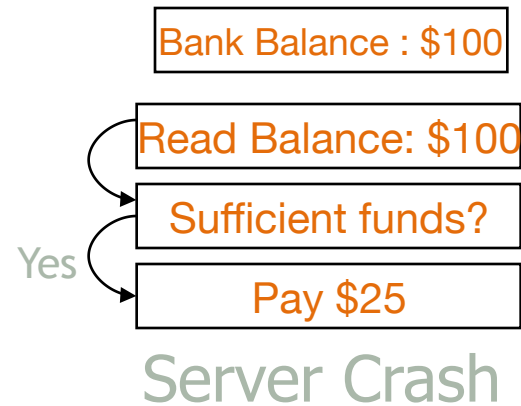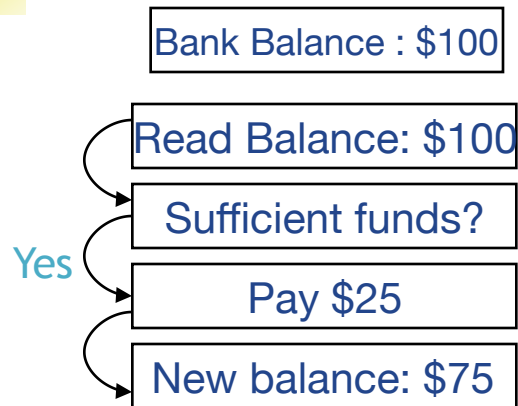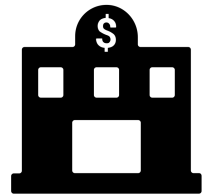
  b) A conceptual schema?

  c) A physical schema?

Think → Pair → Share

# Transactions (I)

| Bank Balance : $100 |
| --- |

| Read Balance: $100 |
| --- |
| Sufficient funds? |
| Pay $25 |
| New balance: $75 |

Yes

| Bank Balance : $100 |
| --- |

| Read Balance: $100 |
| --- |
| Sufficient funds? |
| Pay $25 |

Yes

Server Crash
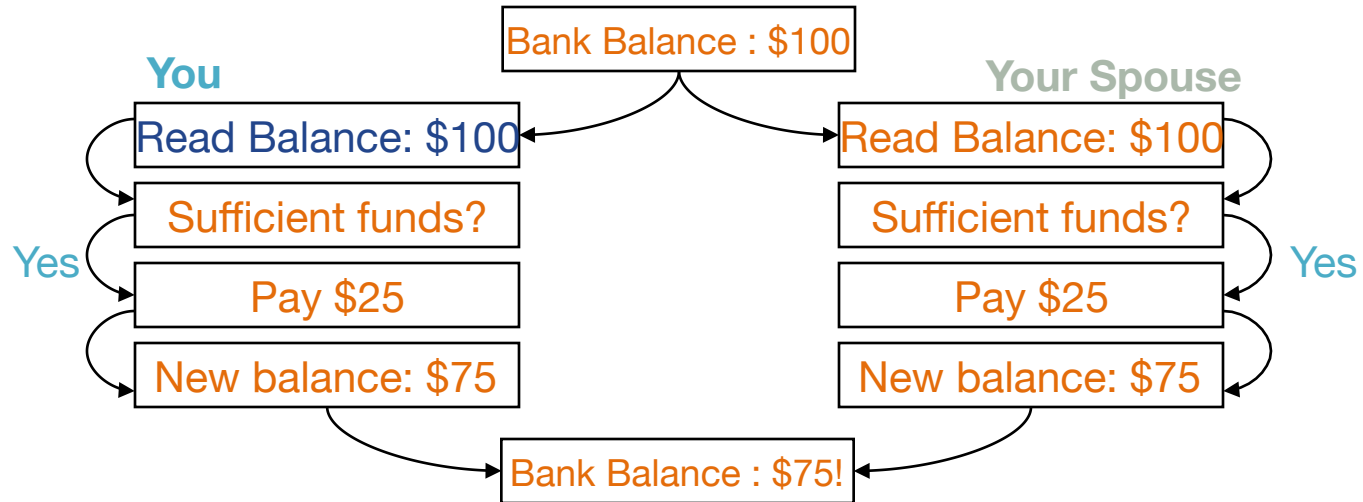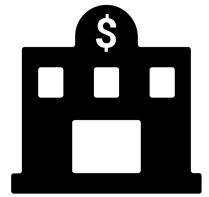
- **Transaction**: any one execution of a user program

  in a DBMS
- Inconsistency caused by incomplete operations
- DBMS ensures atomic operations!
  - i.e., all or nothing!
  - Automatic recovery from crashes!

# Transactions (II)



Bank Balance : $100

**You**

Read Balance: $100

Sufficient funds?

Yes

Pay $25

New balance: $75

**Your Spouse**

Read Balance: $100

Sufficient funds?

Yes

Pay $25

New balance: $75

Bank Balance : $75!

- Inconsistency caused by interleaving actions of different user programs
- DBMS provides the illusion of a "single-user" system
  - Key concept: **Transaction**, an atomic sequence of R/W
  - Concurrency control, transaction management

# Lots of People use DBMS …

- DBMS vendors
- DB application programmers
  - E.g. smart webmasters
- *Database administrator (DBA)*
  - Designs logical /physical schemas
  - Handles security and authorization
  - Data availability, crash recovery
  - Database tuning as needs evolve

**Must understand how a DBMS works!**

# Summary

- DBMS used to maintain, query large datasets.
- Benefits include recovery from system crashes, concurrent access, quick application development, data integrity and security.
- Levels of abstraction give data independence.
- DBAs hold responsible jobs and are well-paid!
- DBMS R&D is one of the most exciting areas in CS.