

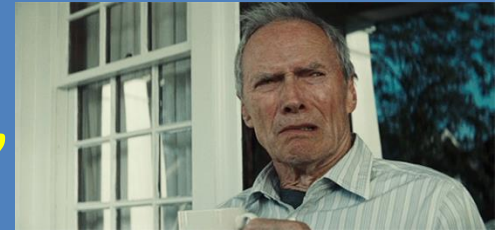
W4111 – Introduction to Databases
Section 003, V03, Fall 2024
Lecture 1: Introduction and Foundational Concepts



W4111 – Introduction to Databases
Section 003, V03, Fall 2024
Lecture 1: Introduction and Foundational Concepts

We will start in a few minutes.

**If you ask me about the waitlist,
I will become very grumpy.**



W4111 Introduction to Databases:

*Faculty do not manage waitlists
for some courses, including W4111.*

*The academic admin staff in the
CS Department manages the waitlist,
priorities and enrollment.*

*You should contact advising email:
ug-advising, ms-advising, or phd-advising
@cs.columbia.edu*

Today's Contents

Contents

- Introduction
 - Logistics, about your instructor, OHs, IAs.
 - Homework, exams.
- ~~Data, Databases, Database Management Systems, Applications~~
- Motivating Examples:
 - Interactive, web application.
 - Data analysis and visualization.
- Database design, Entity-Relationship Model (Part 1)
 - Database design process.
 - ER-Model and diagrams.
- The theory: The *Relational Model* (Part 1)
 - Relational model, schema, keys, schema diagrams.
 - Basics of relational algebra.
- The realization: Structured Query Language (SQL) (Part 1)
 - Basics of Data Definition Language.
 - Basics of Data Manipulation Language (Query).
- Homework 1 – (Initial) Definition and discussion.
- The material is mostly
 - History of databases
 - Motivation for databases
 - Terms
 - Concepts
- Covering in lecture is not a good use of time.
- Just read slides that come with book.
<https://www.db-book.com/slides-dir/index.html>
- Homework assignments and exams will test knowledge and if you have reviewed material.

Introduction

Logistics

Lecture Format, Recitation, Office Hours

- Sessions:
 - Lecture: (In-Person) Lecture: Friday, 10:10AM to 12:40PM (309 Havemeyer).
 - Office hours: Fridays, 8:30AM – 10:00 and 4:00 PM to 5:00 PM (488 CSB)
 - Extra office hours: I hold a lot of extra office hours, usually online, and based on workload around assignment due dates and exams.
 - I have a Zoom meeting and record all lectures and office hours.
- Collaboration/contact:
 - We will use Ed Discussions for collaboration and communication. This is available from the side navigation menu on CourseWorks.
 - The course lectures, sample code, etc. will be in a [GitHub](#) repository.
 - The course [website](#) provides additional information. (In progress)
 - There will be a [course calendar](#) with OHs, assignments, exams, ... (In progress)
- **Note: Your health, safety and well-being are ALWAYS my primary concern. Speak to me if you need special considerations and I will do the best that I can.**

About Your Instructor

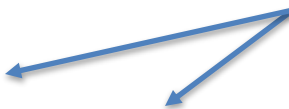
About your Instructor

- 38 years in computer science industry:
 - IBM Fellow.
 - Microsoft Technical Fellow.
 - Chief Technology Officer, CA technologies.
 - Dell Senior Technical Fellow.
 - CTO, Co-Founder, [Seeka.tv](https://seeka.tv).
 - Ansys (current):
 - Ansys Fellow, Chief SW Architect;
 - VP/GM, Cloud, AI, Solutions and Developer Enablement BU (CASEBU)

- Academic experience:

- BA, MS, Ph.D., Computer Science, Columbia University.
- Approx. 18 semesters as an Adjunct Professor.
- Professor of Professional Practice in CS (2018)
- Courses:
 - E1006: Intro. to Computing
 - W4111: Intro. to Databases
 - E6998, E6156: Advanced Topics in SW Engineering (Cloud Computing)

I have taught some version
of this class 10 times.



- Approx. 65 technical publications; Approx. 12 patents.



Personal:

- Two children:
 - College Sophomore.
 - 2019 Barnard Graduate.
- Hobbies:
 - Krav Maga, Black Belt in Kenpo Karate.
 - Former 1LT, [New York Guard](https://www.nyng.org/).
 - Bicycling.
 - Astronomy.
 - Languages:
 - Proficient in Spanish.
 - Learning Arabic.

About the Course

The Course

- From the new/pending Columbia University Bulletin

"Prerequisites: COMS W3134, COMS W3136, or COMS W3137; or the instructor's permission.

The course covers what a database system is, how to design databases effectively and in a principled manner, how to query databases, and how to develop applications using databases: entity-relationship modeling, logical design of relational databases, relational algebra, SQL, database application development, database security, and an overview of query optimization and transaction processing. Additional topics generally include NoSQL, graph, object-relational, and cloud databases, as well as data preparation and cleaning of real-world data. The course offers both programming and non-programming paths for homework and projects, to accommodate students with different programming skills and backgrounds."

- Prerequisites:
 - COMS W3134, COMS W3136, or COMS W3137 are data structures classes. All of these courses require extensive programming, in Java.
 - A course in data structures is helpful for this section of W4111 but not essential. I waive the requirement.
 - We will help you with any data structures knowledge you lack.
- Programming in/for this class:
 - There will be a "non-programming" option, which we will discuss below.
 - For students who want to take the programming track, we will use Python. I will provide motivation for choosing Python below.

Course Objectives

- Have fun, learn a lot and come to appreciate and enjoy some amazing technology. Data and databases have and will change the world.
- Provide a foundation that allows you to succeed in future courses. This is an *introduction* to databases. The technology is crucial for future courses
 - Big data, data analysis, data science
 - Advanced database classes
 - Machine learning
 - Numerical and data analytics in operations research, engineering, economics, finance, life sciences, financial engineering, medicine, etc.
- Enable you to successfully apply the technology in your work and profession.

Have cool stuff to talk about on interviews and in your resumes.

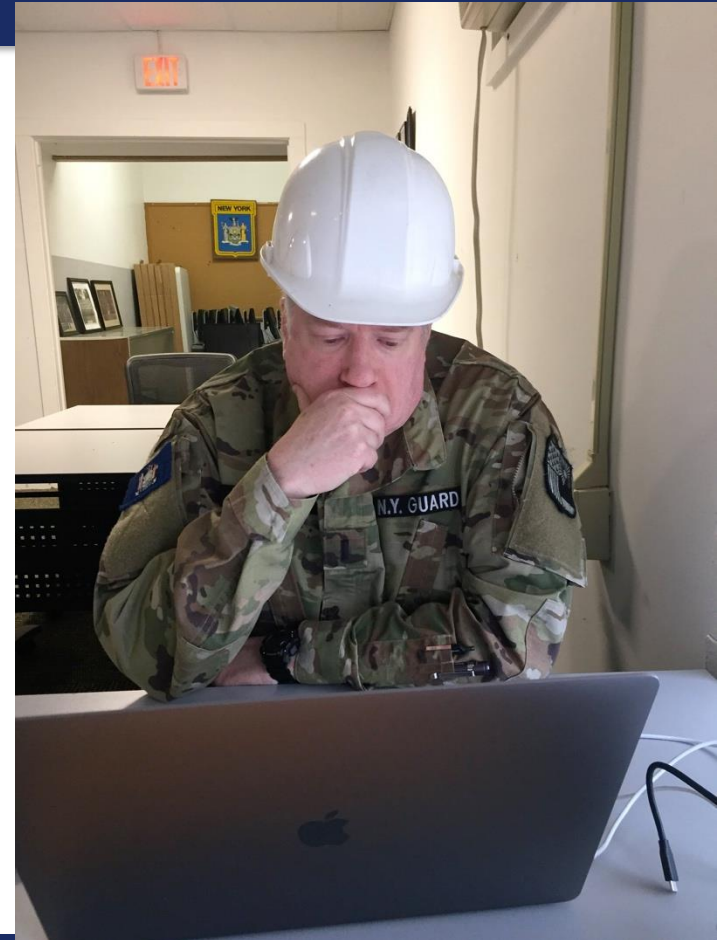
The Course – Value and my Perspective

- This course is foundational, and will teach you the core concepts in
 - Data modeling
 - Data model implementation; Data manipulation.
 - Different database models and database management systems.
 - Implementation and architecture of data centric applications and database management systems.
- ANY non-trivial application
 - Requires a well-designed data model.
 - Implements a data model and manipulates data.
 - Uses a database management system.
- Understanding databases and database management is core to the “hottest fields” in computer science, e.g.
 - Data science
 - Machine learning
 - Intelligent (Autonomous) systems
 - Internet-of-Things
 - Cybersecurity
 - Cloud Computing
- Database, database application, etc. skills are increasingly central to many disciplines, including:
 - Economics.
 - IEOR, financial engineering.
 - Mechanical and electrical engineering.
 - Medicine, pharmaceuticals
 - ...
- University courses on databases sometimes focus on theory and abstract concepts. This course will cover theory, but there will be an increased emphasis on:
 - Practical, hands-on applications of databases.
 - Patterns and best practices.
 - Developing and understanding database centric applications.
 - Using databases and various tools for data analysis, visualization and insight.
- **Personal perspective**
 - A large percent of my career has been spent figuring out or leading teams that figured out how to model, implement and manipulate data.
 - I have used the information in this class more than anything else I have learned.
 - This will likely be true for you.

Surprising Example

Example: This is a photo of me using a database during COVID-19 mobilization.

- Match service members
- To JTF-HQ requests for personnel
- Based on assignment needs
- And service member
 - Skills
 - Availability
 - Tracked in a DBMS.
- The hardhat is because DB usage can be very dangerous 🧐.
- No joking: I had to build an application that used:
 - Relational DBMS.
 - Python, Jupyter Notebook.
 - Google Sheets with Apps Script.
 - Google Forms.
- This course's technology and these skills are surprisingly applicable.



Modules

Each section of W4111 is slightly different based on student interest and professor's focus. There is a common, core syllabus. Professors cover topics in different orders and grouping based on teaching style.

This section of W4111 has four modules:

- **Foundational concepts (50% of semester):** This module covers concepts like data models, relational model, relational databases and applications, schema, normalization, ... The module focuses on the relational model and relational databases. The concepts are critical and foundational for all types of databases and data centric applications.
- **Database management system architecture and implementation (10%):** This module covers the software architecture, algorithms and implementation techniques that allow [databases management systems](#) to deliver functions. Topics include memory hierarchy, storage systems, caching/buffer pools, indexes, query processing, query optimization, transaction processing, isolation and concurrency control.
- **NoSQL – “Not Only SQL” databases (20%):** This module provides motivation for [“NoSQL”](#) data models and databases, and covers examples and use cases. The module also includes cloud databases and databases-as-a-service.
- **Data Enabled Decision Support (20%):** This module covers data warehouses, data import and cleanse, OLAP, Pivot Tables, Star Schema, reporting and visualization, and provides an overview of analysis techniques, e.g. clustering, classification, analysis, mining.

Environment

Course Resources and Development Environment

- Recommended textbook:
 - *Database System Concepts. Seventh Edition.* (ISBN 9780078022159)
 - There is a website associated with the textbook: <https://www.db-book.com/>. The site has:
 - Slides for each chapter.
 - Example data.
 - Textbooks are expensive. You can easily get through the course using website, lecture material,
- Install a new, most recent, isolated/single user instance of Anaconda environment.
 - **Install just for yourself and within your home directory.**
 - You must install the most [recent version](#) for Python 3.
 - You can isolate the new instance from other instances to avoid conflicts, or set up custom environments.
- Development environments: Students are entitled to a free, annual [JetBrains professional license](#).
 - Students have had problems with MySQL Workbench. We will use [DataGrip](#). Please install.
 - Installing [PyCharm](#) is recommended for all and required for the programming track.
- Install [MySQL Server Community Edition](#).
 - When prompted, choose [legacy authentication](#) method.
 - Set your root password to **dbuserdbuser**. (**Remember your user name and password**)

Homework Exams Grading

Assignments, Exams, Grading

- To be updated.

To Program or Not To Program, ...

- From the department's guidance for the course:
 - "To accommodate the diverse backgrounds of the students who take this class, all sections of the class should include a non-programming option for projects and assignments. We (NOTE: Does not include me) have successfully offered such an option in some sections of the class for many years: students can either program a web application to interface with a DB of their own design, or alternatively follow the non-programming option and come up with a quantifiably more detailed DB design/data.
- What does not constitute programming?
 - Writing queries for a relational, graph, document, or other data management system does not constitute programming and can be expected of all students.
- What constitutes programming?
 - Reading more than a handful (1-5) lines of code, Writing Python code that is not directly required to write a query."
- Why I previously only did a programming track:
 - Any non-trivial use of data and databases in any field requires programming. This is not Star Trek. You cannot shout, "Computer! Analyze ..."
 - You might think, "But, I do not want to program. I want to go into financial services, private equity, medicine, ... I can just use tools like Pandas or Tableau". You will do some programming in these jobs for complex scenarios. Get over it.
- Tracks: Programming and non-programming.
 - There is a common core that involves understanding concepts, modeling data, using databases, etc.
 - Exams are the same for both tracks. Mix of written questions and practical exercises.
 - Homework assignments: There is a common core on all assignments for both tracks. Additionally,
 - Programming track will incrementally build a database centric, cloud-based web application.
 - Non-programming track will do a more complex, database centric project.

Homework Assignments

- My homework assignments are:
 - Open ended.
 - Vaguely specified.
- You will complain. I will listen sympathetically.
 - You will savage me on the instructor reviews and CULPA.
 - My professional colleagues and I will laugh at you behind your back.
- Management, clients, partners, etc.
 - Do not understand technology as well as we do.
 - You will get requirements like, "I want it colored mauve because that has more RAM."
- Converting vague requests into a useful, meaningful project that we can implement is what we do. Get over it.

You are senior people at a top school. You need to define concrete solutions and approaches from ambiguity. No one showed Thomas Edison a light bulb and then said, "Build this."

Donald Ferguson on 2019-01-29

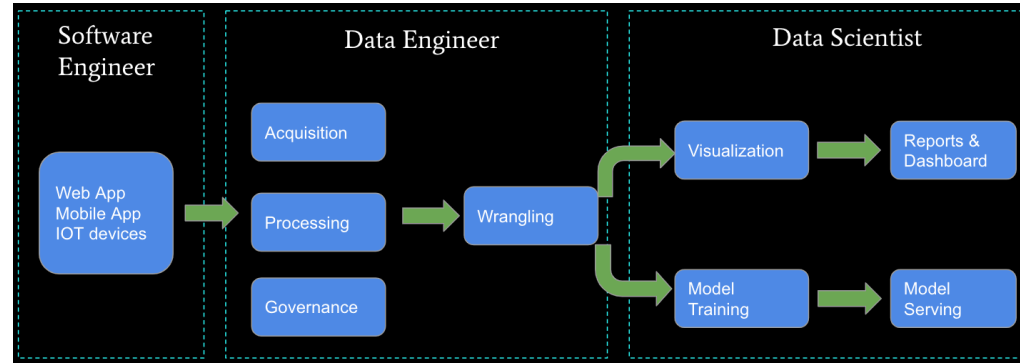
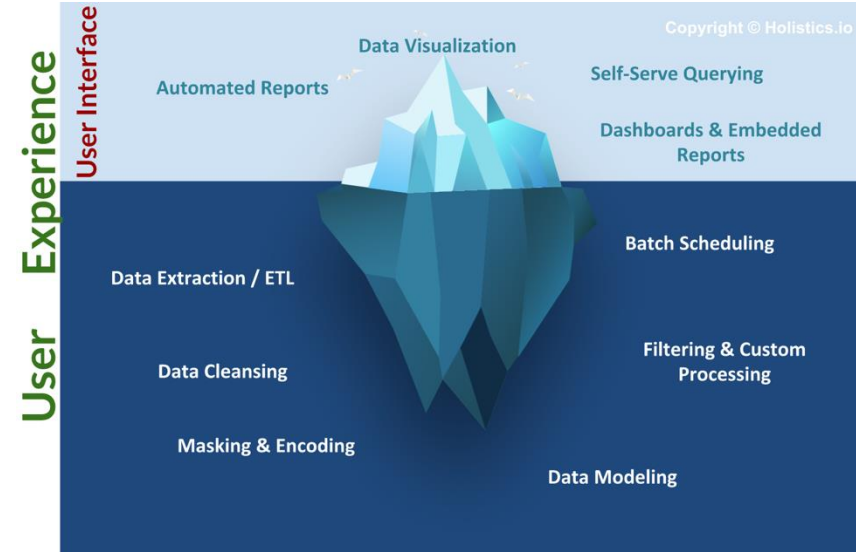


Database Applications

Two Common Database Applications

- Operational/Interactive:
 - Users and roles can create, retrieve, update, search and delete “records.”
 - Examples: SSOL, ATMs,
- Business Intelligence, Decision Support, ...:
 - Users can perform complex queries and analyze a lot of data to generate a report, make a decision,
 - Examples: Build AI/ML training data, dashboards,
- Some of our major datasets this semester will be:
 - IMDB: <https://developer.imdb.com/non-commercial-datasets/>
 - Game of Thrones: <https://developer.imdb.com/non-commercial-datasets/>
 - Lahman’s Baseball Dataset: <http://seanlahman.com/>
 -
- We will build a simple web application and do some data engineering.

Business Intelligence, Insight, Analysis,

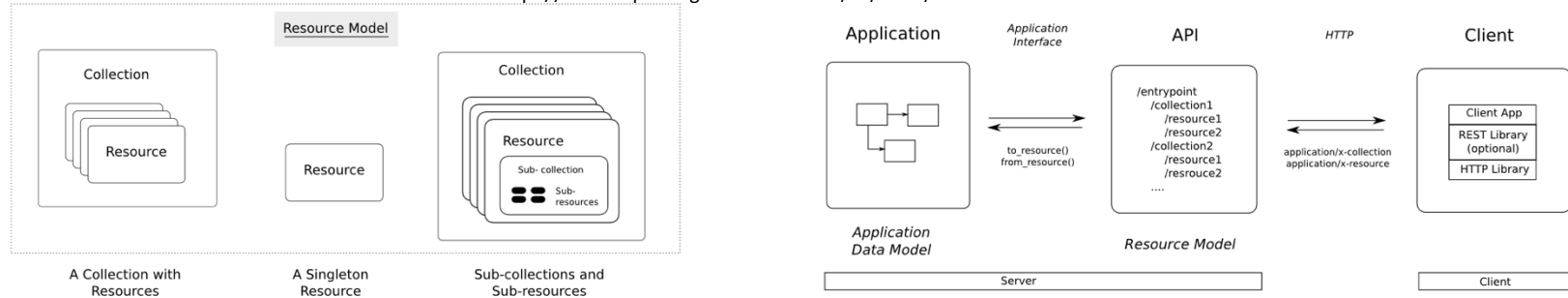


- The “fun” stuff in data science and AI/ML is the “tip of the iceberg.”
- Data engineering is a necessary condition for producing analyzable data. This is often more than 80% of the hard work.
- We will do some small data engineering projects in this course.

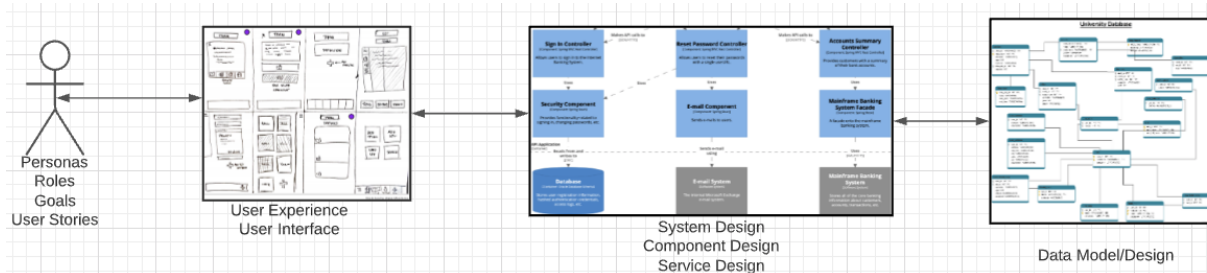
Web Application Problem Statement

- We must build a system that supports create, retrieve, update and delete for IMDB and Game of Thrones Datasets.
- This requires implementing *create, retrieve, update and delete (CRUD)* for resources.

<https://restful-api-design.readthedocs.io/en/latest/resources.html>



- We will design, develop, test and deploy the system iteratively and continuously.
- There are four core domains.

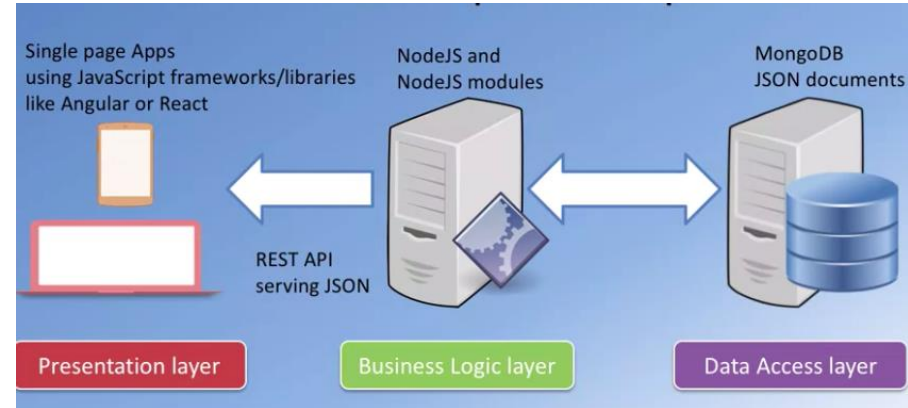


- In this course,
- We focus on the data dimension.
- We will get some insight into the other dimensions.

Interactive/Operational

- “A full stack web developer is a person who can develop both client and server software. In addition to mastering HTML and CSS, he/she also knows how to:
 - Program a browser (like using JavaScript, jQuery, Angular, or Vue)
 - Program a server (like using PHP, ASP, Python, or Node)
 - Program a database (like using SQL, SQLite, or MongoDB)”

https://www.w3schools.com/whatis/whatis_fullstack.asp
- We will do a simple full stack app.
 - Three databases:
 - MySQL
 - MongoDB
 - Neo4j
 - The application tier will be Python and FastAPI.
 - The web UI will be Angular.
 - The primary focus is the data layer and application layer that access it.
 - I will provide a simple UI and template.



Approach for First Few Lectures

Lectures and Topics

► Chapter 1 Introduction

▼ PART ONE RELATIONAL LANGUAGES

► Chapter 2 Introduction to the Relational M...

► Chapter 3 Introduction to SQL

► Chapter 4 Intermediate SQL

► Chapter 5 Advanced SQL

▼ PART TWO DATABASE DESIGN

► Chapter 6 Database Design Using the E-R ...

► Chapter 7 Relational Database Design

▼ PART THREE APPLICATION DESIGN AND DEV...

► Chapter 8 Complex Data Types

- Chapter 1:
 - Is interesting, but something easily learned from reading slides, docs, etc.
 - **You are responsible for reading the slides and independent study.**
- Books and standard syllabuses tend to be sequential. Cover and complete topics one at a time.
- My view is that there is a core conceptual model with three realizations:
 - Entity, relationships,; the concepts.
 - ER design and modeling.
 - Implementation:
 - Relational model/algebra.
 - SQL.
 - Resource/REST oriented.
 -
- My approach is to incrementally and iteratively:
 - Introduce concepts.
 - Explain the realizations.
 - Because implementing a real system requires the approach of concept, design, implement in several layers.

Database Modeling, ER Modeling (I)



Design Phases

- Initial phase -- characterize fully the data needs of the prospective database users.
- Second phase -- choosing a data model
 - Applying the concepts of the chosen data model
 - Translating these requirements into a conceptual schema of the database.
 - A fully developed conceptual schema indicates the functional requirements of the enterprise.
 - Describe the kinds of operations (or transactions) that will be performed on the data.

DFF Comments:

- We you see slides with this formatting, the come directly from the presentations associated with the textbook. (<https://www.db-book.com/db7/slides-dir/index.html>)
- The number at the bottom is of the form chapter.slide_no.
- I try to put my comments, modifications and annotations in red text, or inside a red rectangle/callout.



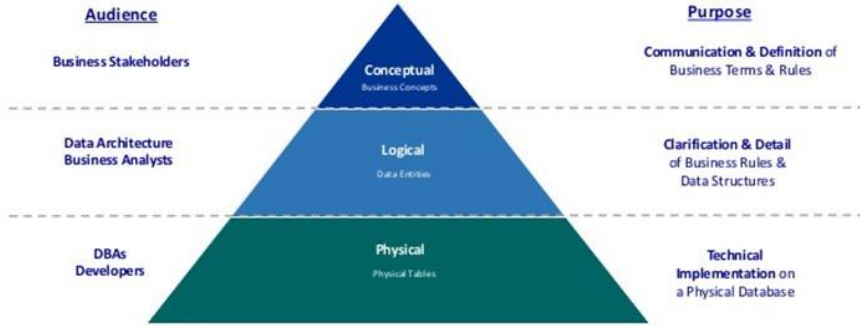
Design Phases (Cont.)

- Final Phase -- Moving from an abstract data model to the implementation 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

A Common and my Approach: Conceptual → Logical → Physical

<https://ehikioya.com/conceptual-logical-physical-database-modeling/>

Levels of Data Modeling



- It is easy to get carried away with modeling. You can spend all your time modeling and not actually build the schema.
- We will use the approaches in class.
- Mostly to understand concepts and patterns.

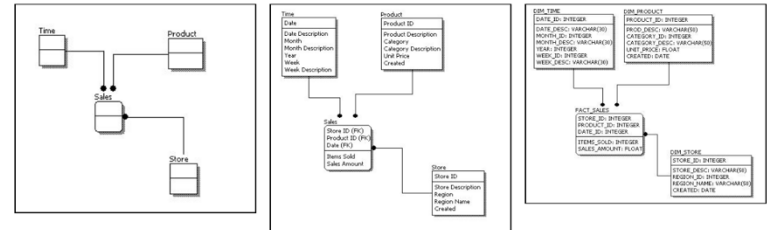
<https://www.1keydata.com/datawarehousing/data-modeling-levels.html>

Feature	Conceptual	Logical	Physical
Entity Names	✓	✓	
Entity Relationships	✓	✓	
Attributes		✓	
Primary Keys		✓	✓
Foreign Keys		✓	✓
Table Names			✓
Column Names			✓
Column Data Types			✓

Conceptual Model Design

Logical Model Design

Physical Model Design



<https://www.1keydata.com/datawarehousing/data-modeling-levels.html>



ER model -- Database Modeling

- The ER data model was developed to facilitate database design by allowing specification of an **enterprise schema** that represents the overall logical structure of a database.
- The ER data model employs three basic concepts:
 - entity sets,
 - relationship sets,
 - attributes.
- The ER model also has an associated diagrammatic representation, the **ER diagram**, which can express the overall logical structure of a database graphically.



Entity Sets

COMS W4111 002 01 2024

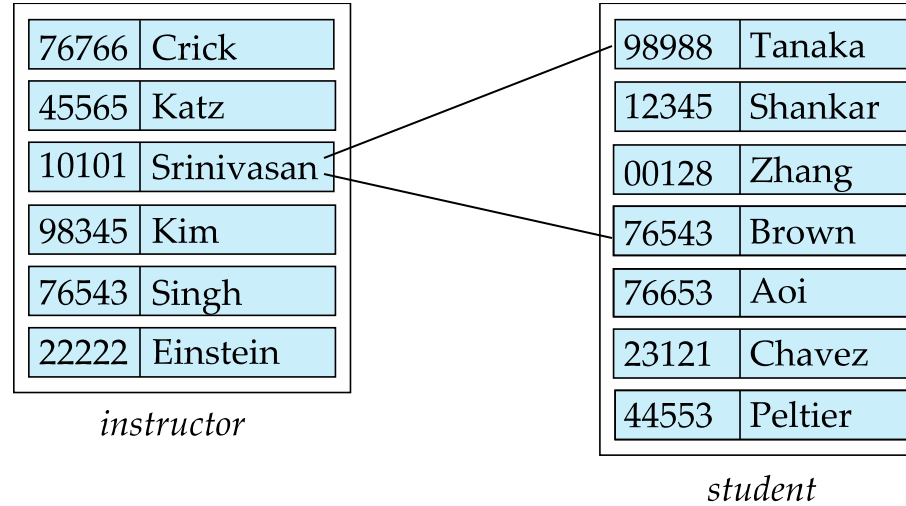
- An **entity** is an object that exists and is distinguishable from other objects.
 - Example: specific person, company, event, plant
- An **entity set** is a set of entities of the **same type** that share the same properties.
 - Example: set of all persons, companies, trees, holidays
- An entity is represented by a set of attributes; i.e., descriptive properties **possessed by all** members of an entity set.
 - Example:
$$\text{instructor} = (ID, name, salary)$$
$$\text{course} = (course_id, title, credits)$$
- A subset of the attributes form a **primary key** of the entity set; i.e., uniquely identifying each member of the set.

DFF Comments:

- Some of these statements apply primarily to OO systems and the relational/SQL models.
- A motivation for “No SQL” is to relax the constraints.



Entity Sets -- *instructor* and *student*



(10101, 98988)

(10101, 76543)



Relationship Sets

- A **relationship** is an association among several entities

Example:

44553 (Peltier)	<u>advisor</u>	22222 (<u>Einstein</u>)
<i>student</i> entity	relationship set	<i>instructor</i> entity

- A **relationship set** is a mathematical relation among $n \geq 2$ entities, each taken from entity sets

$$\{(e_1, e_2, \dots, e_n) \mid e_1 \in E_1, e_2 \in E_2, \dots, e_n \in E_n\}$$

where (e_1, e_2, \dots, e_n) is a relationship

- Example:

$$(44553, 22222) \in \text{advisor}$$

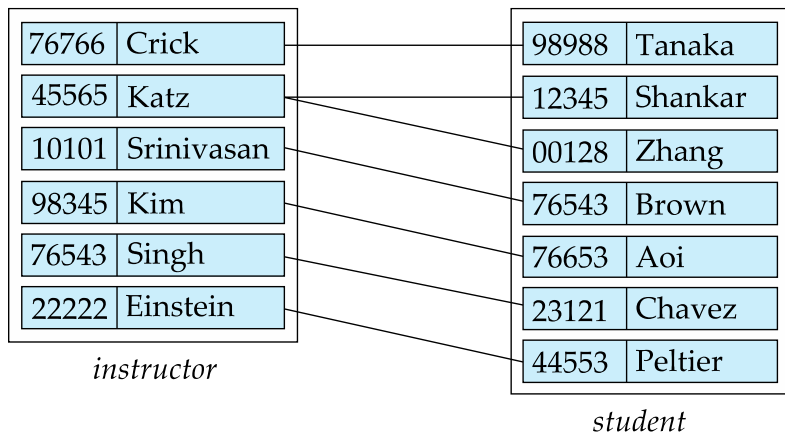
DFF Comments:

- Nobody thinks about relationships this way.
- There is no idea so simple that a DB professor cannot make it confusing, usually by using math.



Relationship Sets (Cont.)

- Example: we define the relationship set *advisor* to denote the associations between students and the instructors who act as their advisors.
- Pictorially, we draw a line between related entities.



DFF Comment:

- In this diagram, the lines are the relationship set.
- Many DB models use a 3rd entity set to represent complex relationships.

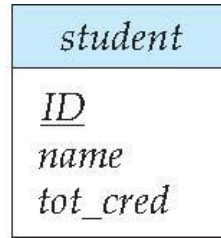
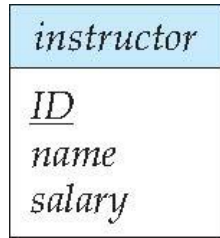
DFF Comments:

- Nobody draws the diagrams this way, but ...
- Sometimes thinking this way helps understand other ways to depict the concept.

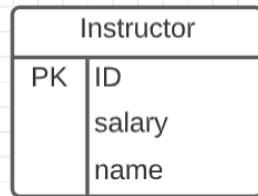


Representing Entity sets in ER Diagram

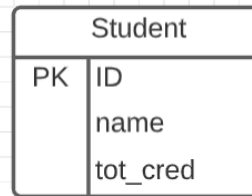
- Entity sets can be represented graphically as follows:
 - Rectangles represent entity sets.
 - Attributes listed inside entity rectangle
 - Underline indicates primary key attributes



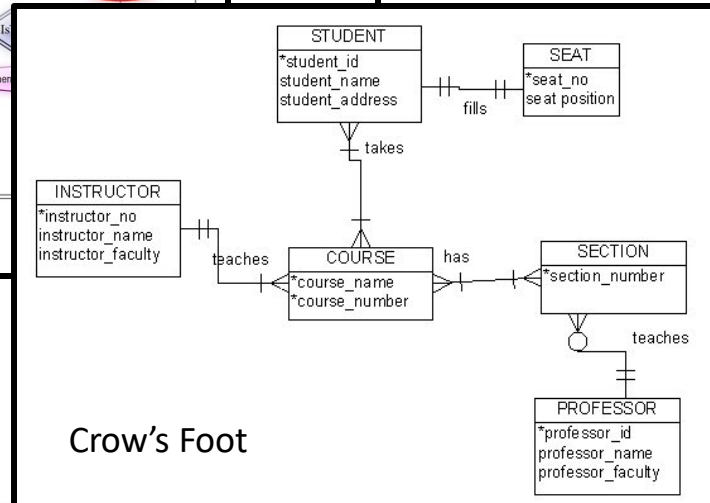
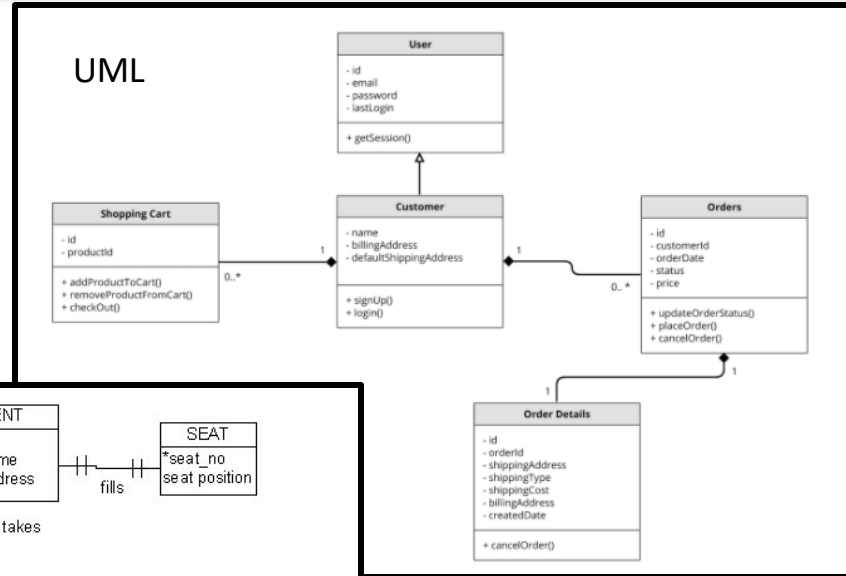
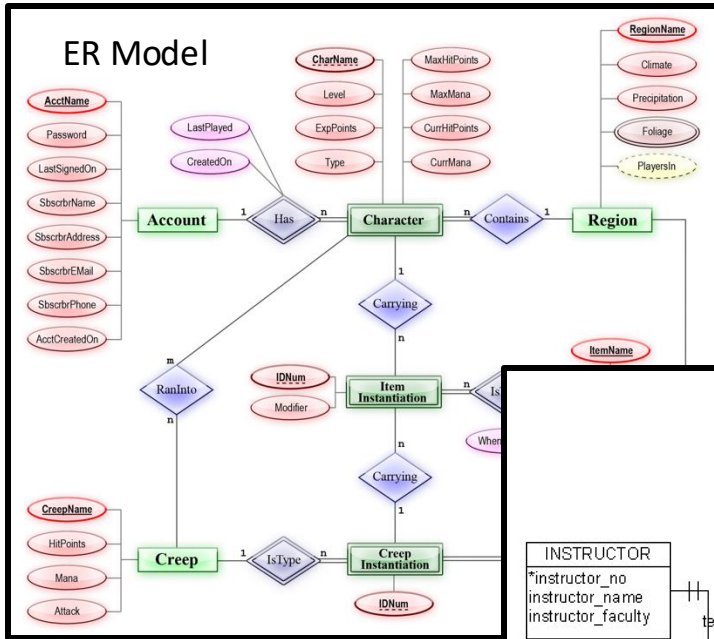
“Visual modeling is the use of semantically rich, graphical and textual design notations to capture software designs. A notation, such as UML, allows the level of abstraction to be raised, while maintaining rigorous syntax and semantics. In this way, it improves communication in the design team, as the design is formed and reviewed, allowing the reader to reason about the design, and it provides an unambiguous basis for implementation.”



Crow's
Foot
Notation



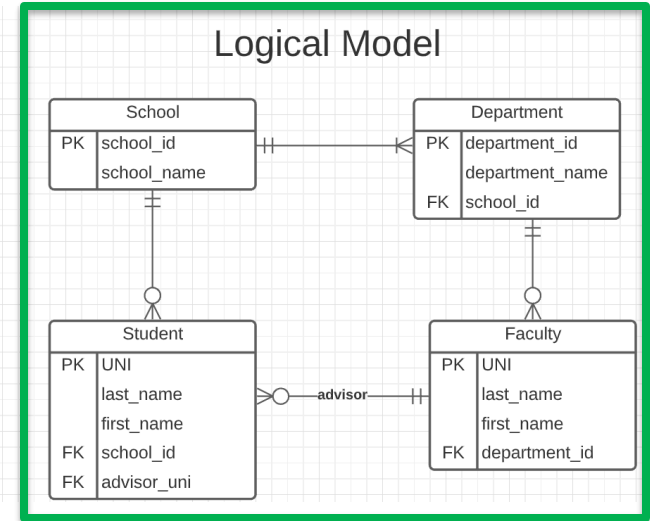
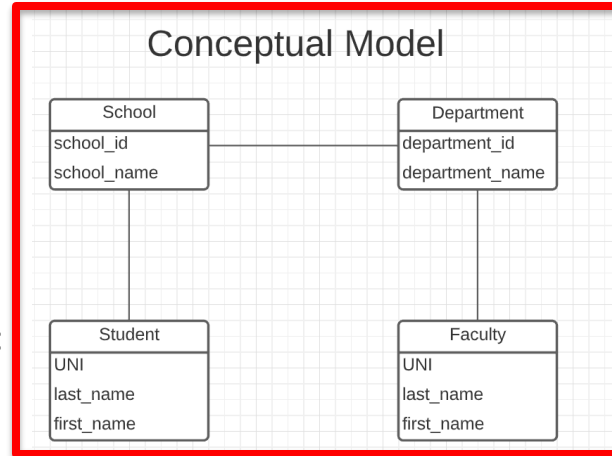
Visual Notation – Many Notations



- “Other,” i.e. PowerPoint is the most common modeling notation.
- It is easy to get “carried away.”
- The trick is to do “just enough modeling.”
- I mostly use Crow’s Foot
 - It is “just enough”
 - But lacks some capabilities.
- The book uses ER notation.

Do a Slightly More complex University Database ER Model

- The model has:
 - Four entity sets:
 - *School*
 - *Department*
 - *Student*
 - *Faculty*
 - Three relationship sets:
 - *Student-School*
 - *School-Department*
 - *Faculty-Department*.

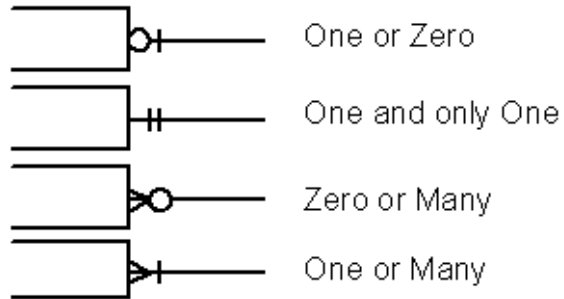


- This is the level of detail I want when I ask for:
 - Conceptual Model diagram.
 - Logical Model diagram.
- Some online tools with “free,” constrained usage.
 - Lucidchart (<https://www.lucidchart.com/>)
 - Vertabelo (<https://vertabelo.com/>)

Notation has Precise Meaning

- Attribute annotations:
 - PK = Primary Key
 - FK = Foreign Key
- Line annotations:
 - We will spend a lot of time discussing keys.
 - We will start in a couple of slides.

Summary of Crow's Foot Notation



- We will learn over time and there are good tutorials (<https://www.lucidchart.com/pages/er-diagrams>) to help study and refresh.

What Does this Mean? Let's Get Started

Primary Key means that the value occurs **at most once**.

This is a statement about the domain, not the actual data currently in the table.

School Code	School Name
CC	Columbia College
SEAS	Fu Foundation School of Engineering and Applied Science
GSAS	Graduate School of Arts and Sciences
GS	General Studies
...	...

Foreign Key means that if a value occurs in `school_id` for any row, there must be a row in School with that key.

UNI	Last name	First name	school_id
dff9	Ferguson	Donald	CC
js11	Smith	John	GS
jp9	Public	James	CC
bb101	Baggins	Bilbo	CC
...

The line notations mean:

- A student is related to EXACTLY ONE school.
- A School may be related to 0, 1 or many students.



ER model -- Database Modeling

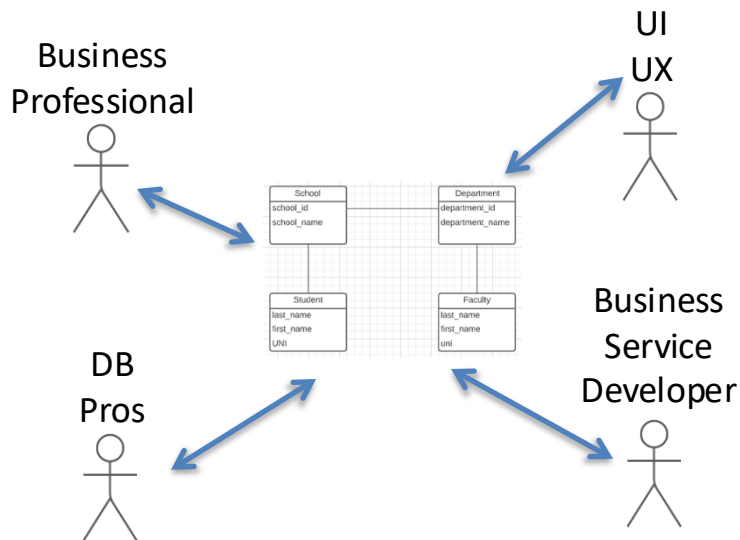
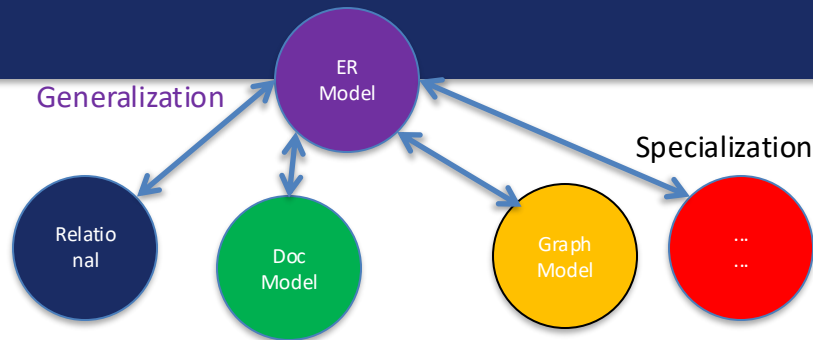
- The ER data model was developed to facilitate database design by allowing specification of an **enterprise schema** that represents the overall logical structure of a database.
- The ER data model employs three basic concepts:
 - entity sets,
 - relationship sets,
 - attributes.
- The ER model also has an associated diagrammatic representation, the **ER diagram**, which can express the overall logical structure of a database graphically.

DFF Comments:

- The book and slides do not do a great job of motivating the ER model or ER diagrams.
- Why do people and teams think about or use the ER model and modeling?

ER Model and ER Modeling

- ER Model: Agility, Separation of Concerns
 - ER model is a generalization that most DB models implement in some form.
 - Using the ER model enables:
 - Thinking about and collaborating on design with getting bogged down in details.
 - Enable flexible choices about how to realize/Implement data.
- ER Diagrams: Communication, Quality, Precision
 - With a little experience, everyone can understand and ER diagram.
 - Easier to discuss and collaborate on application's data than showing SQL table definitions, JSON,
 - People think visually. That is why we have whiteboards. ER diagrams are precise and unambiguous.
 - Guides you to think about relationships, keys, ... And prevents “re-dos” later in the process. It is easier to fix a diagram than a database schema.



ER Modeling – Reasonably Good Summary

Advantages of ER Model

Conceptually it is very simple: ER model is very simple because if we know relationship between entities and attributes, then we can easily draw an ER diagram.

Better visual representation: ER model is a diagrammatic representation of any logical structure of database. By seeing ER diagram, we can easily understand relationship among entities and relationship.

Effective communication tool: It is an effective communication tool for database designer.

Highly integrated with relational model: ER model can be easily converted into relational model by simply converting ER model into tables.

Easy conversion to any data model: ER model can be easily converted into another data model like hierarchical data model, network data model and so on.

Disadvantages of ER Model

Limited constraints and specification

Loss of information content: Some information be lost or hidden in ER model

Limited relationship representation: ER model represents limited relationship as compared to another data models like relational model etc.

No representation of data manipulation: It is difficult to show data manipulation in ER model.

Popular for high level design: ER model is very popular for designing high level design

No industry standard for notation

<https://pctechanicalpro.blogspot.com/2017/04/advantages-disadvantages-er-model-dbms.html>

Note:

- If you get to use Google to help with take home exams, HW, etc.
- I get to use Google to help with slides.

Theory: Relational Model (I)



Relational Model

- All the data is stored in various tables.
- Example of tabular data in the relational model



Ted Codd
Turing Award 1981

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

- The “relation” is the “table.”
 - In my big space of pieces of data. *ID, name, dept_name, salary* are somehow related.
 - This causes confusion, because the ER and other models use “relation” to mean something else.
- Core concepts:
 - Relation
 - Tuple (Row)
 - Column (Attribute)



Example of a *Instructor* Relation

00123
123

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

tuples
(or rows)

Integer
-213
(integers, >0)



Attribute

- The set of allowed values for each attribute is called the **domain** of the attribute
- Attribute values are (normally) required to be **atomic**; that is, indivisible
- The special value **null** is a member of every domain. Indicated that the value is “unknown”
- The null value causes complications in the definition of many operations

Ferguson, Donald

Ferguson Donald

COMSW4111

DFF Comments:

- I will explain the importance of atomic attributes and null in examples.
- Atomic and use of Null is important!



Relations are Unordered

- Order of tuples is irrelevant (tuples may be stored in an arbitrary order)
- Example: *instructor* relation with unordered tuples

<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



Database Schema

- Database schema -- is the logical structure of the database.
- Database instance -- is a snapshot of the data in the database at a given instant in time.
- Example:
 - schema: *instructor* (*ID*, *name*, *dept_name*, *salary*)
 - Instance:

<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



Keys

- Let $K \subseteq R$
- K is a **superkey** of R if values for K are sufficient to identify a unique tuple of each possible relation $r(R)$
 - Example: $\{ID\}$ and $\{ID, name\}$ are both superkeys of *instructor*.
- Superkey K is a **candidate key** if K is minimal
Example: $\{ID\}$ is a candidate key for *Instructor*
- One of the candidate keys is selected to be the **primary key**.
 - which one?
- **Foreign key** constraint: Value in one relation must appear in another
 - **Referencing** relation
 - **Referenced** relation
 - Example: *dept_name* in *instructor* is a foreign key from *instructor* referencing *department*

Notation

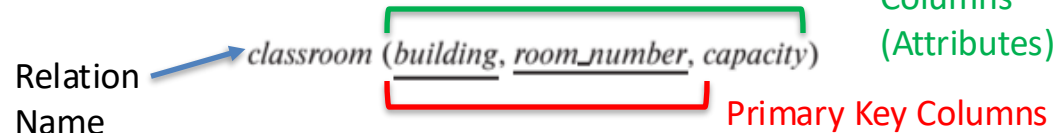
Classroom relation

building	room_number	capacity
Packard	101	500
Painter	100	125
Painter	514	10
Taylor	3128	70
Watson	100	30
Watson	120	50

classroom schema

It is customary to list the primary key attributes of a relation schema before the other attributes; for example, the *dept_name* attribute of *department* is listed first, since it is the primary key. Primary key attributes are also underlined.

Consider the *classroom* relation:



- The primary key is a *composite key*. Neither column is a key (unique) by itself.
- Keys are statements about all possible, valid tuples and not just the ones in the relation.
 - Capacity is unique in this specific data, but clearly not unique for all possible data.
 - In this domain, there cannot be two classrooms with the same building and room number.
- Relation schema:
 - Underline indicates a primary key column. There is no standard way to indicate other types of key.
 - We will use **bold** to indicate foreign keys.
 - You will sometimes see things like *classroom*(*building:string*, *room_number:number*, *capacity:number*)

Observations

- Keys:
 - Will be baffling. It takes time and experience to understand/appreciate.
 - There are many, many types of keys with formal definitions.
 - I explain the formality but focus on the concepts and applications.
- No one uses the formal, relational model. So, why do we study it?
 - Is very helpful when understanding concepts that we cover later in the course, especially query optimization and processing.
 - There are many realizations of the model and algebra, and understanding the foundation helps with understanding language/engine capabilities.
 - The model has helped with innovating new approaches, and you may innovate in data and query models in your future.



Relational Query Languages

- Procedural versus non-procedural, or declarative
- “Pure” languages:
 - Relational algebra
 - Tuple relational calculus
 - Domain relational calculus
- The above 3 pure languages are equivalent in computing power
- We will concentrate in this chapter on relational algebra
 - Not turning-machine equivalent
 - Consists of 6 basic operations

DFF Comments:

- You will sometimes see other operator, e.g. \leftarrow Assignment.
- Relational algebra focuses on *retrieve*. You can sort of do Create, Update, Delete.
- The SQL Language, which we will see, extends relational algebra.



Relational Algebra

- A procedural language consisting of a set of operations that take one or two relations as input and produce a new relation as their result.
- Six basic operators
 - select: σ
 - project: Π
 - union: \cup
 - set difference: $-$
 - Cartesian product: \times
 - rename: ρ



Select Operation

- The **select** operation selects tuples that satisfy a given predicate.
- Notation: $\sigma_p(r)$
- p is called the **selection predicate**
- Example: select those tuples of the *instructor* relation where the instructor is in the “Physics” department.
 - Query

$\sigma_{dept_name=“Physics”}(instructor)$

- Result

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
22222	Einstein	Physics	95000
33456	Gold	Physics	87000



Select Operation (Cont.)

- We allow comparisons using

$=, \neq, >, \geq, <, \leq$

in the selection predicate.

- We can combine several predicates into a larger predicate by using the connectives:

\wedge (**and**), \vee (**or**), \neg (**not**)

- Example: Find the instructors in Physics with a salary greater \$90,000, we write:

$\sigma_{dept_name='Physics' \wedge salary > 90,000} (instructor)$

- Then select predicate may include comparisons between two attributes.
 - Example, find all departments whose name is the same as their building name:
 - $\sigma_{dept_name=building} (department)$



Project Operation

- A unary operation that returns its argument relation, with certain attributes left out.
- Notation:

$$\Pi_{A_1, A_2, A_3, \dots, A_k}(r)$$

where A_1, A_2 are attribute names and r is a relation name.

- The result is defined as the relation of k columns obtained by erasing the columns that are not listed
- Duplicate rows removed from result, since relations are sets



Project Operation (Cont.)

- Example: eliminate the *dept_name* attribute of *instructor*
- Query:

$\Pi_{ID, name, salary} (instructor)$

- Result:

<i>ID</i>	<i>name</i>	<i>salary</i>
10101	Srinivasan	65000
12121	Wu	90000
15151	Mozart	40000
22222	Einstein	95000
32343	El Said	60000
33456	Gold	87000
45565	Katz	75000
58583	Califieri	62000
76543	Singh	80000
76766	Crick	72000
83821	Brandt	92000
98345	Kim	80000



Composition of Relational Operations

- The result of a relational-algebra operation is relation and therefore of relational-algebra operations can be composed together into a **relational-algebra expression**.
- Consider the query -- Find the names of all instructors in the Physics department.

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

- Instead of giving the name of a relation as the argument of the projection operation, we give an expression that evaluates to a relation.

The Dreaded Relax Calculator

- Let's look at an online tool that you will use.
- RelaX (<https://dbis-uibk.github.io/relax/calc/local/uibk/local/0>)
- The calculator:
 - Has an older version of the data from the recommended textbook.
(<https://dbis-uibk.github.io/relax/calc/gist/4f7866c17624ca9dfa85ed2482078be8/relax-silberschatz-english.txt/0>)
 - You can also upload new data.
- Some queries:
 - $\sigma \text{ dept_name} = \text{'Comp. Sci.'} \vee \text{dept_name} = \text{'History'}$ (department)
 - $\pi \text{ name, dept_name}$ (instructor)
 - $\pi \text{ ID, name}$ ($\sigma \text{ dept_name} = \text{'Comp. Sci.'}$ (instructor))

Structured Query Language (I)



History

- IBM Sequel language developed as part of System R project at the IBM San Jose Research Laboratory
- Renamed Structured Query Language (SQL)
- ANSI and ISO standard SQL:
 - SQL-86
 - SQL-89
 - SQL-92
 - SQL:1999 (language name became Y2K compliant!)
 - SQL:2003
- Commercial systems offer most, if not all, SQL-92 features, plus varying feature sets from later standards and special proprietary features.
 - Not all examples here may work on your particular system.



SQL Parts

- DML -- provides the ability to query information from the database and to insert tuples into, delete tuples from, and modify tuples in the database.
- integrity – the DDL includes commands for specifying integrity constraints.
- View definition -- The DDL includes commands for defining views.
- Transaction control –includes commands for specifying the beginning and ending of transactions.
- Embedded SQL and dynamic SQL -- define how SQL statements can be embedded within general-purpose programming languages.
- Authorization – includes commands for specifying access rights to relations and views.

DDL

DML

SQL Language Statements

The core SQL language statements are:

- SELECT: Implements both σ , π
 - INSERT
 - UPDATE
 - DELETE
 - CREATE TABLE
 - ALTER TABLE
 - JOIN, which is an operator within SELECT.
- Many, if not most, SQL statements:
 - Implement multiple relational algebra expressions.
 - Cannot easily (or at all) be represented in relational algebra.

```
 $\pi$  ID, name (  
     $\sigma$  dept_name='Comp. Sci.' (instructor)  
)  
=  
SELECT ID, name FROM instructor  
WHERE  
dept_name='Comp. Sci.'
```



Basic Query Structure

- A typical SQL query has the form:

```
select  $A_1, A_2, \dots, A_n$   
from  $r_1, r_2, \dots, r_m$   
where  $P$ 
```

- A_i represents an attribute
 - R_i represents a relation
 - P is a predicate.
- The result of an SQL query is a relation.

Note:

- The SELECT ... FROM ... WHERE ... Combines two relational operators, σ and Π .
- Actually, it also combines other operators, e.g. \times



The select Clause

- The **select** clause lists the attributes desired in the result of a query
 - corresponds to the projection operation of the relational algebra
- Example: find the names of all instructors:
select *name*
from *instructor*
- NOTE: SQL names are case insensitive (i.e., you may use upper- or lower-case letters.)
 - E.g., *Name* \equiv *NAME* \equiv *name*
 - Some people use upper case wherever we use bold font.



The select Clause (Cont.)

- SQL allows duplicates in relations as well as in query results.
- To force the elimination of duplicates, insert the keyword **distinct** after select.
- Find the department names of all instructors, and remove duplicates

```
select distinct dept_name  
from instructor
```

- The keyword **all** specifies that duplicates should not be removed.

```
select all dept_name  
from instructor
```



The select Clause (Cont.)

- An asterisk in the select clause denotes “all attributes”

```
select *  
from instructor
```

- An attribute can be a literal with no **from** clause

```
select '437'
```

- Results is a table with one column and a single row with value “437”
- Can give the column a name using:

```
select '437' as FOO
```

- An attribute can be a literal with **from** clause

```
select 'A'  
from instructor
```

- Result is a table with one column and N rows (number of tuples in the *instructors* table), each row with value “A”



The select Clause (Cont.)

- The **select** clause can contain arithmetic expressions involving the operation, +, −, *, and /, and operating on constants or attributes of tuples.
 - The query:

```
select ID, name, salary/12  
from instructor
```

would return a relation that is the same as the *instructor* relation, except that the value of the attribute *salary* is divided by 12.

- Can rename “*salary/12*” using the **as** clause:

```
select ID, name, salary/12 as monthly_salary
```




The where Clause

- The **where** clause specifies conditions that the result must satisfy
 - Corresponds to the selection predicate of the relational algebra.
- To find all instructors in Comp. Sci. dept

```
select name  
from instructor  
where dept_name = 'Comp. Sci.'
```

- SQL allows the use of the logical connectives **and**, **or**, and **not**
- The operands of the logical connectives can be expressions involving the comparison operators **<**, **<=**, **>**, **>=**, **=**, and **<>**.
- Comparisons can be applied to results of arithmetic expressions
- To find all instructors in Comp. Sci. dept with salary > 80000

```
select name  
from instructor  
where dept_name = 'Comp. Sci.' and salary > 80000
```



The from Clause

- The **from** clause lists the relations involved in the query
 - Corresponds to the Cartesian product operation of the relational algebra.
- Find the Cartesian product *instructor X teaches*

select *
from *instructor, teaches*

- generates every possible instructor – teaches pair, with all attributes from both relations.
 - For common attributes (e.g., *ID*), the attributes in the resulting table are renamed using the relation name (e.g., *instructor.ID*)
- Cartesian product not very useful directly, but useful combined with where-clause condition (selection operation in relational algebra).



Examples

- Find the names of all instructors who have taught some course and the course_id
 - **select** *name, course_id*
from *instructor , teaches*
where *instructor.ID = teaches.ID*

- Find the names of all instructors in the Art department who have taught some course and the course_id
 - **select** *name, course_id*
from *instructor , teaches*
where *instructor.ID = teaches.ID and instructor.dept_name = 'Art'*



Data Definition Language

The SQL data-definition language (DDL) allows the specification of information about relations, including:

- The schema for each relation.
- The type of values associated with each attribute.
- The Integrity constraints
- The set of indices to be maintained for each relation.
- Security and authorization information for each relation.
- The physical storage structure of each relation on disk.



Create Table Construct

- An SQL relation is defined using the **create table** command:

create table *r*

($A_1 D_1, A_2 D_2, \dots, A_n D_n,$
(integrity-constraint₁),
...,
(integrity-constraint_k))

- *r* is the name of the relation
 - each A_i is an attribute name in the schema of relation *r*
 - D_i is the data type of values in the domain of attribute A_i
- Example:

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

Some Worked Example

Switch to Demos: Jupyter Notebook and DataGrip

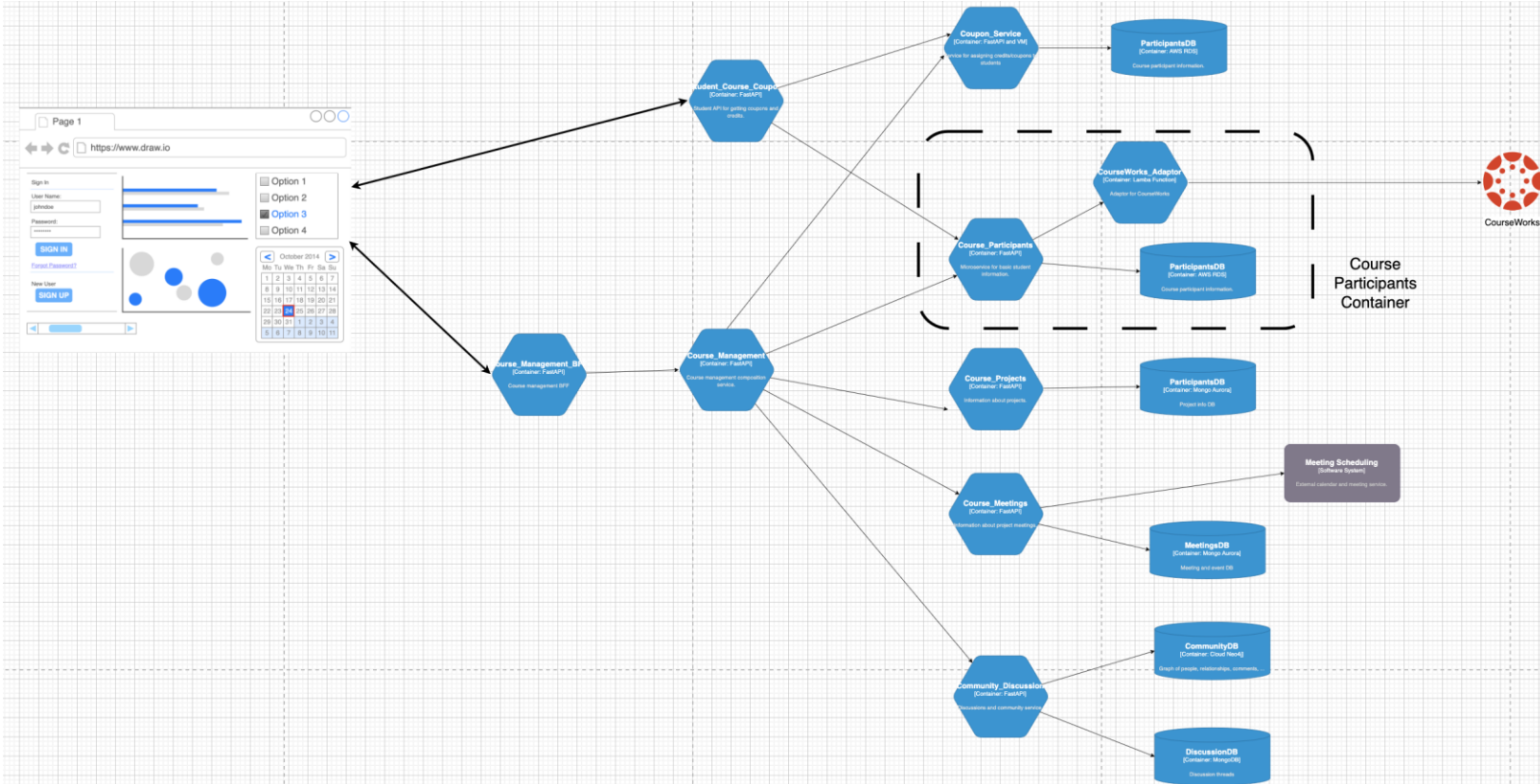
- Drawing a simple data model in Lucidchart.
- Generating sample data in Mockaroo.
- Loading CSV data using DataGrip.
- Editing and modifying schema in SQL.
- Running some queries.
- Doing some relational algebra.

Simple First Scenario

Scenario

- I am building an application to help me manage my classes, especially the cloud computing courses.
- The primary complexities are:
 - Assigning and managing project teams and students on teams.
 - Tracking and commenting on project documents.
 - Meeting management.
 -
- This is primarily an interaction application using multiple databases.
- There are some data engineering and analysis problems, however. For example, loading data from CourseWorks.
- CourseWorks has an API. I can write code to call the API and extract the information.

Microservice Application



CourseWorks Canvas API

Canvas LMS - REST API and Extensions Documentation

Expand all

› Basics

› OAuth2

▼ Resources

All Resources
Reference

Access Tokens

Account Calendars

Account Domain
Lookups

Account Notifications

Account Reports

Accounts

Accounts (LTI)

Admins

Analytics

Announcement
External Feeds

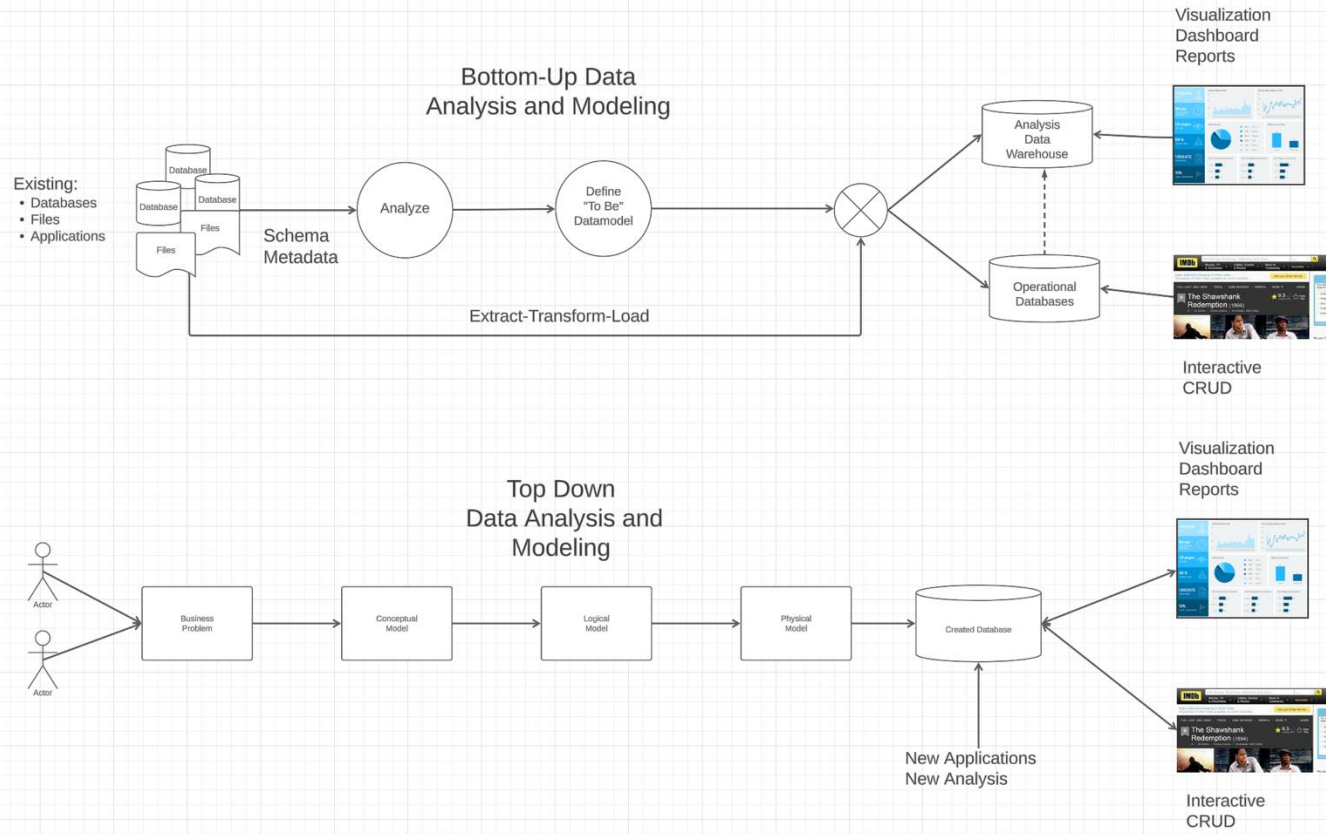
Announcements

API Token Scopes

All API Resources

Access Tokens	Show an access token Create an access token Update an access token Delete an access token
Account Calendars	List available account calendars Get a single account calendar Update a calendar Update several calendars List all account calendars Count of all visible account calendars
Account Domain Lookups	Search account domains
Account Notifications	Index of active global notification for the user Show a global notification Close notification for user Create a global notification Update a global notification
Account Reports	List Available Reports Start a Report Index of Reports Status of a Report Delete a Report
Accounts	List accounts Get accounts that admins can manage Get accounts that users can create courses in List accounts for course admins Get a single account Settings List environment settings Permissions Get the sub-accounts of an account Get the Terms of Service Get help links Get the manually-created courses sub-account for the domain root account List active courses in an account Update an account Delete a user from the root account Restore a deleted user from a root account
Subaccounts	Create a new sub-account Delete a sub-account
Accounts (LTI)	Get account
Admins	Make an account admin Remove account admin List account admins List my admin roles
Analytics	Get department-level participation data Get department-level grade data Get department-level statistics Get department-level statistics, broken down by subaccount Get course-level participation data Get course-level assignment data Get course-level student summary data Get user-in-a-course-level participation data Get user-in-a-course-level assignment data Get user-in-a-course-level messaging data

Data Modeling and Engineering



Data “Engineering” Pipeline

- The CourseWorks API:
 - Returns a big, complex JSON document for each Course-Section
 - Returns a bigger, more complex JSON document containing students in Course-Section.
- I need to transform the data:
 - Remove uninteresting information for my application.
 - Define to be relational model.
 - Load the data.
- Show some of the notebook functions and underlying Python scripts.

Interactive Demo

- `/Users/donald.ferguson/Dropbox/000/00-Current-Repos/student-course-application/courseworks-adaptor-service/adaptor/courseworks_data_processing.ipynb`
- Interactive application:
 - `/Users/donald.ferguson/Dropbox/000/00-Current-Repos/student-course-application/course-section-example`
 - `/Users/donald.ferguson/Dropbox/000/00-Current-Repos/student-course-application/dff_framework`

End of Lecture

HW0 and HW1 will come out over the weekend.