

# *W4111 – Introduction to Databases*

## *Section 002, Fall 2021*

*Lecture 1: Introduction and Foundational Concepts (V1.0)*



# *W4111 – Introduction to Databases*

## *Section 002, Fall 2021*

*Lecture 1: Introduction and Foundational Concepts (V1.0)*

*We will start in a couple of minutes*



# *Today's Contents*

# Contents

- Introduction
    - Logistics, about your instructor, OHs, TAs.
    - Changes from prior semesters based on student feedback.
    - Homework, exams.
  - Introduction
    - Introductory concepts: data, databases, database management systems.
    - Examples: Gives you a feel for programming and non-programming projects, applications, application architectures.
  - 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: Structure Query Language (SQL) (Part 1)
    - Basics of Data Definition Language.
    - Basics of Data Manipulation Language (Query).
  - Homework 0 – 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/db7/slides-dir/index.html>

*Introduction*  
*Logistics*  
*Changes Based on Feedback*  
*Exams, Homework*

# *Logistics*

# Waitlist

- The administrative staff of the Dept. of Computer Science **directly** manages enrollment in the W4111 sections and the wait lists. This is true of several courses.
  - COMS W4111 is **ALWAYS** oversubscribed and has huge wait lists.
  - The course is a requirement for several majors and tracks.
  - There are **complex rules** about prioritizing students for course enrollment, and also which students can take the hybrid section.
- Faculty follow department policies and do not manage waitlists, despite the list that identifies the waitlist as “Instructor Managed.”
- Contacts: You must work with administrative leaders on enrollment:  
Computer Science: Cynthia Meekins, [cynthia@cs.columbia.edu](mailto:cynthia@cs.columbia.edu))

# Lecture Format, Recitation, Office Hours

- Sessions:
  - Lecture: (In-Person) Lecture: Friday, 10:10 to 12:40
  - Recitation: Periodic, optional recitation on Saturday, 10:00 to 11:30.
    - Either in-person/online or just online. Will announce in advance.
    - Usually held when there is student interest, homework due, exams, etc.
- Office hours:
  - In-person: Friday: 0830 to 1000, 4:00 PM – 6:00 PM; and by-request as needed.
  - Extra office hours: I hold a lot of extra office hours, usually based on workload around assignment due dates and exams.
    - I will be on campus when possible for students that want to speak with me or need more focused help.
    - Many of the extra OHs will be online only.
- Collaboration/contact:
  - We will use [Ed Discussions](#), which is available from CourseWorks. This replaces Piazza.
  - I am usually monitoring [Slack](#). Subscribe to the channel #w4111-f21.
  - The course lectures, sample code, etc. will be in a [GitHub](#) repository.
  - The course website provides additional information.
- **Note: Your health, safety and well-being are ALWAYS my primary concern. These are bizarre, stressful times. Speak to me if you need special considerations and I will do the best that I can.**

# *About Your Instructor*

# About your Instructor

- 35 years in computer science industry:
  - IBM Fellow
  - Microsoft Technical Fellow.
  - Chief Technology Officer, CA technologies.
  - Dell Senior Technical Fellow.
  - CTO, Co-Founder, [Seeka.tv](#).
  - Ansys Fellow
- Academic experience:
  - BA, MS, Ph.D., Computer Science, Columbia University.
  - Approx. 15 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)
- Approx. 65 technical publications; Approx. 12 patents.



## Personal:

- Two children:
  - College freshman (Ithaca College).
  - 2019 Barnard Graduate.
- Hobbies:
  - Krav Maga, Black Belt in Kenpo Karate.
  - 1LT, [New York Guard](#).
  - 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.

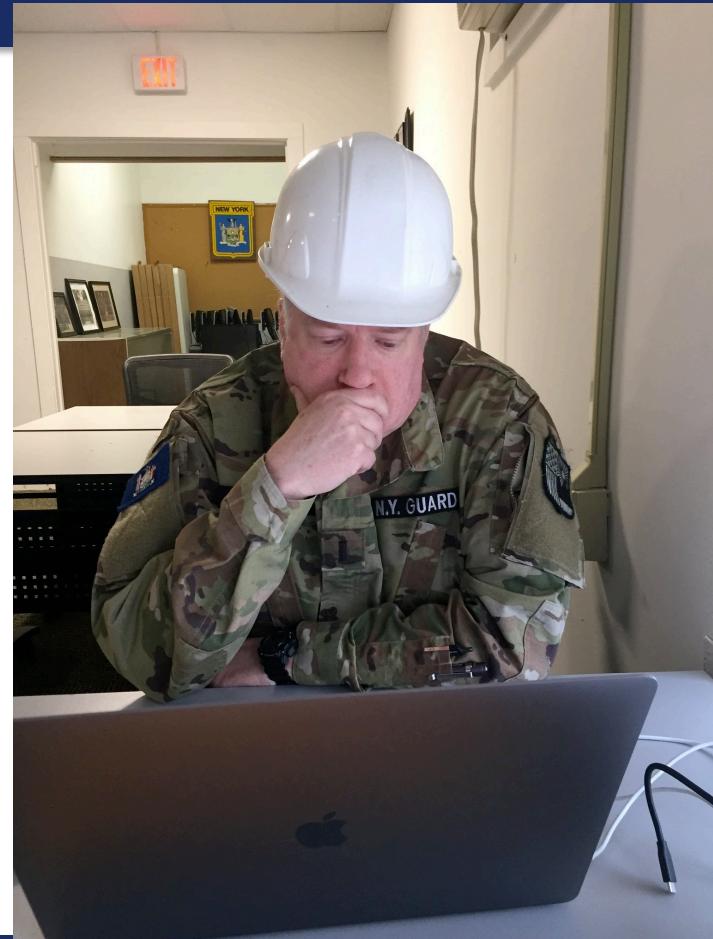
# The Course – Value and my Perspective

- 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 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/db7/>
  - 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.
  - 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 try [DataGrip](#). Please install.
  - Installing [PyCharm](#) is recommended for all and required for the programming track.
- Install [MySQL Service Community Edition](#).
  - When prompted, choose [legacy authentication](#) method.
  - Set your root password to **dbuserdbuser**

# *Changes from Previous Semesters*

# Observations from Last Semester

- Student evaluation:
  - I pay a lot of attention to students' evaluations and comments.
  - Fall 2020 evaluations:
    - Course overall quality was 3.49 (out of 5). My past average has consistently been over 4.1.
    - Instructor overall quality was 3.66. My past average has consistently been over 4.1.
  - Some recurring (negative) themes:
    - "This is the most disorganized course I've taken."
    - HW assignments are vague and unclear.
    - Too much work.
- We will try to make improvements this semester.
  - Much, much better adherence to a HW and exam cadence/schedule.
  - We will make the HW assignments definitions clearer, but not as clear as you want.
  - Relying very heavily on the head TAs for homework/exam definition, explanation, ... ....
  - Will have periodic checkpoint polls to get feedback and change suggestions.

*Homework  
Exams  
Grading*

# Assignments, Exams, Grading

- Point value of assignments and exams
  - 50%: Homework assignments:
    - 4 HWs, approximately one HW every two-three weeks.
    - Each is worth either 10 or 15 points based on difficulty.
  - Homework format:
    - Common to all tracks:
      - Questions requiring written answers and diagrams.
      - Implement/execute/test various database operations required to solve a use case/question.
      - Format will be an iPython/Jupyter Notebook.
    - Track specific: Incremental development of a project:
      - Programming track – web application.
      - Non-programming track – data analysis/visualization.
- All tracks will learn core concepts and get experience with other tracks' focus topics.
- Exams: Both the midterm and final have an in-class and take-home element.
  - 20% of grade is midterm exam score.
  - 30% is final exam score.
- 97-100 points is an A. There are opportunities to get extra-credit.

# 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 previous 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.
- **BUT, I WILL SIGNIFICANTLY IMPROVE CLARITY AND DETAILS.**

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

Omonbude Emmanuel  
@BUDESCODE

To replace programmers with Robots, clients will have to accurately describe what they want.

We're safe.

7:25 PM · 20 Jul 20 · Twitter for Android

# *Examples*

# *IMDB*

## *(Programming Track Example)*

# Most of Us are Familiar with IMDB

The screenshot shows the IMDB mobile website. At the top, there's a search bar with a yellow border and a magnifying glass icon. To the right of the search bar is a "Sign In" button. Below the search bar is a banner for USAA Auto Insurance with the text "MUST-SEE SAVINGS: UP TO 30% OFF WITH USAA SAFEPILOT". The main content area features a scene from Game of Thrones with two characters in a dimly lit setting. A large yellow arrow points from the search bar on the left towards the search results on the right. On the left side of the main content, there's a "WHAT TO WATCH" section with a play button icon and the text "Type ‘Tom Han’ in the search field." Below this is a video thumbnail titled "8 Most Surprising Moments From ‘Game of Thrones’ Episodes That Deserve a Rewatch" with a duration of 3:39. At the bottom left, there's a "Browse trailers >" link.

The screenshot shows the IMDB search results for "Tom Han". At the top, there's a banner for USAA Auto Insurance with the text "SAFE PILOTS SAVE UP TO 30% ON AUTO INSURANCE". To the right of the banner is a movie poster for "Ringers" with the text "SAVINGS". Below the banner, the search results are displayed under the heading "Results for 'Tom Han'". There are two sections: "Names" and "Titles". Under "Names", there are two entries: "Tom Hanks (Producer, Cast Away (2000))" and "Aamir Khan (I) (Actor, Like Stars on Earth (2007)) nickname "Tom Hanks of India"". There is a link "View: More name matches". Under "Titles", there are two entries: "Killing Lincoln (2013) (TV Movie) aka "Tom Hanks: Die Lincoln-Verschwörung"" and "Tom Hanks: A League of His Own (2016)". There is a link "View: More title matches". On the far right, there are vertical navigation links: "Cat", "Searc", "categ", "All", and "Name".

# Downloadable IMDB Data

**title.akas.tsv.gz** - Contains the following information for titles:

- titleId (string) - a tconst, an alphanumeric unique identifier of the title
- ordering (integer) - a number to uniquely identify rows for a given titleId
- title (string) - the localized title
- region (string) - the region for this version of the title
- language (string) - the language of the title
- types (array) - Enumerated set of attributes for this alternative title. One or more of the following: "alternative", "dvd", "festival", "tv", "video", "working", "original", "imdbDisplay". New values may be added in the future without warning
- attributes (array) - Additional terms to describe this alternative title, not enumerated
- isOriginalTitle (boolean) - 0: not original title; 1: original title

**title.basics.tsv.gz** - Contains the following information for titles:

- tconst (string) - alphanumeric unique identifier of the title
- titleType (string) - the type/format of the title (e.g. movie, short, tvseries, tvepisode, video, etc)
- primaryTitle (string) - the more popular title / the title used by the filmmakers on promotional materials at the point of release
- originalTitle (string) - original title, in the original language
- isAdult (boolean) - 0: non-adult title; 1: adult title
- startYear (YYYY) - represents the release year of a title. In the case of TV Series, it is the series start year
- endYear (YYYY) - TV Series end year. '\N' for all other title types
- runtimeMinutes - primary runtime of the title, in minutes
- genres (string array) - includes up to three genres associated with the title

**title.crew.tsv.gz** - Contains the director and writer information for all the titles in IMDb.

Fields include:

- tconst (string) - alphanumeric unique identifier of the title
- directors (array of nconsts) - director(s) of the given title
- writers (array of nconsts) - writer(s) of the given title

**title.episode.tsv.gz** - Contains the tv episode information. Fields include:

- tconst (string) - alphanumeric identifier of episode
- parentTconst (string) - alphanumeric identifier of the parent TV Series
- seasonNumber (integer) - season number the episode belongs to
- episodeNumber (integer) - episode number of the tconst in the TV series

**title.principals.tsv.gz** - Contains the principal cast/crew for titles

- tconst (string) - alphanumeric unique identifier of the title
- ordering (integer) - a number to uniquely identify rows for a given titleId
- nconst (string) - alphanumeric unique identifier of the name/person
- category (string) - the category of job that person was in
- job (string) - the specific job title if applicable, else '\N'
- characters (string) - the name of the character played if applicable, else '\N'

**title.ratings.tsv.gz** - Contains the IMDb rating and votes information for titles

- tconst (string) - alphanumeric unique identifier of the title
- averageRating - weighted average of all the individual user ratings
- numVotes - number of votes the title has received

**name.basics.tsv.gz** - Contains the following information for names:

- nconst (string) - alphanumeric unique identifier of the name/person
- primaryName (string) - name by which the person is most often credited
- birthYear - in YYYY format
- deathYear - in YYYY format if applicable, else '\N'
- primaryProfession (array of strings) - the top-3 professions of the person
- knownForTitles (array of tconsts) - titles the person is known for

<https://www.imdb.com/interfaces/>

# Conceptual Model – Application Architecture

<https://dzone.com/articles/architectural-shift-in-web-applications-with-emerg>

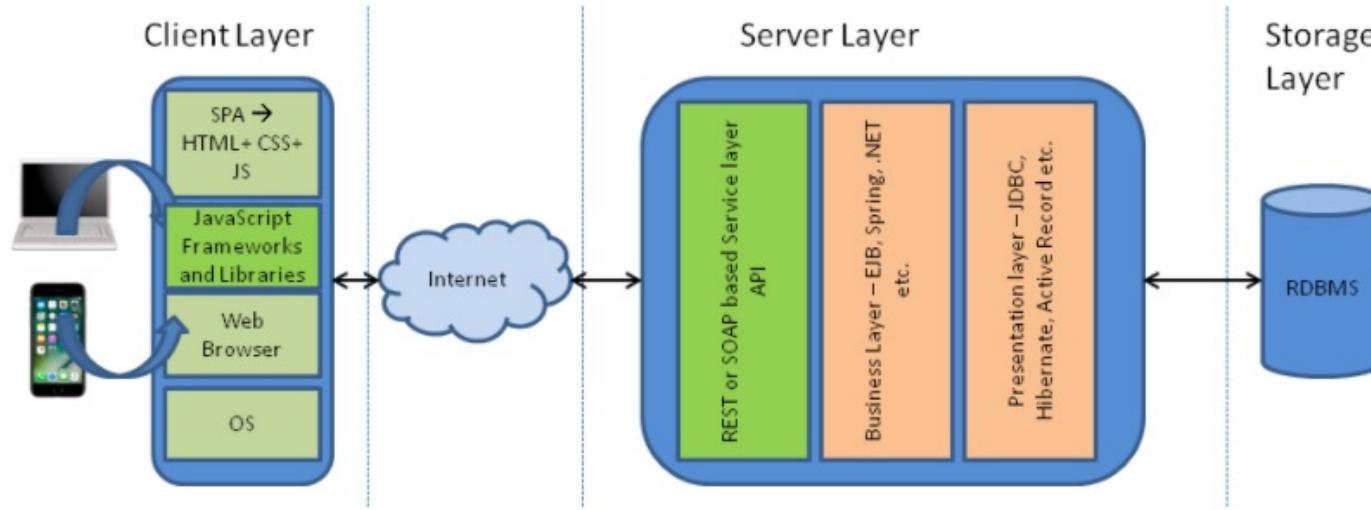
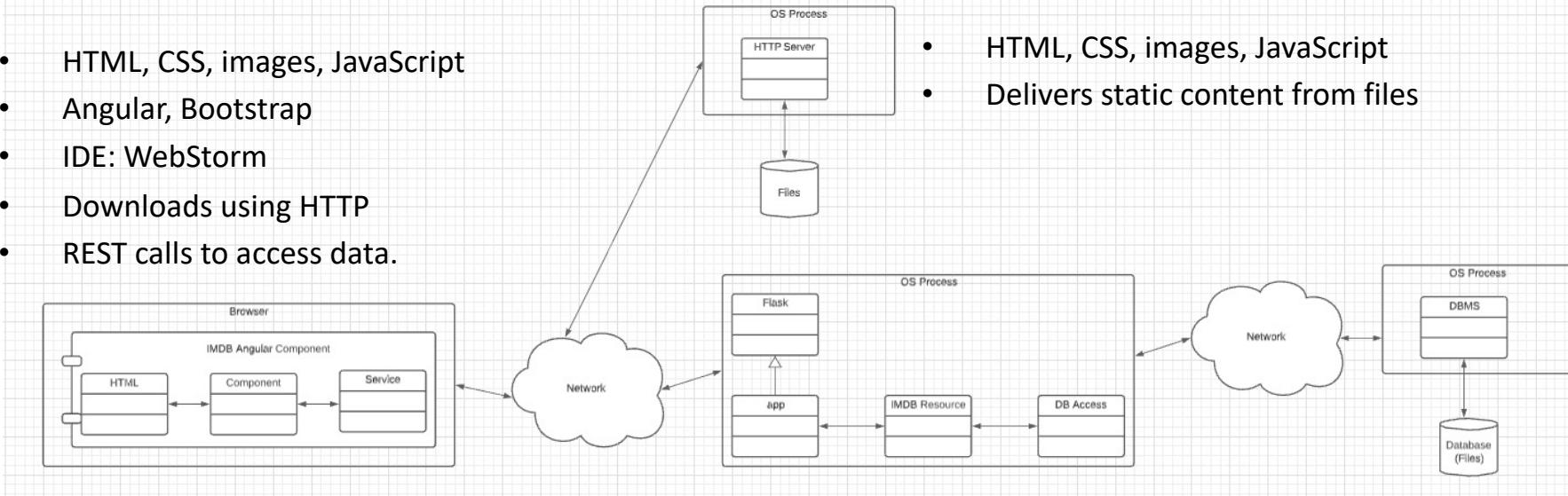


Diagram 2: The moving of the Web Layer from the Server to the Client

- Browser client application with HTML, CSS, JavaScript, ... ...
- Middle-tier server application:
  - Receives API calls from client application.
  - Implements application logic and makes database calls.
  - Returns result.
- Database Management System receives DB calls and returns data.
- This is a “full stack (web) application.”
- Part of the programming track HW and projects is to build a full stack app.
- Starting from a template I provide to simplify non-DB tasks.

# Explanation of Example Code

- HTML, CSS, images, JavaScript
- Angular, Bootstrap
- IDE: WebStorm
- Downloads using HTTP
- REST calls to access data.



- Flask application server framework
- Python application code
- IDE: PyCharm
- Interacts with database over network connection using SQL statements.
- MySQL Community Server
- Implements SQL operations
- Maintains data in files

# Switch to Demo

## A Simple IMDB

- Single page browser application.
- Web application (server).
- Database schema and query.

## Show some code.

- Browser
- Python, Flask
- SQL database
- REST requests

### IMDB Information

Some information about [IMDB](#) artists will go here! Type in an artist's name:

Tom Hank

Primary

Primary Name	Birth Year	Death Year	Primary Professions	Known For Titles
Tom Hanks	1956		producer,actor,soundtrack	tt0109830,tt0094737,tt0162222,tt0120815
Tom Hankason			actor,actor	tt0072759
Tom Hankins			actor	tt11066358
Tom Hanke			sound_department,composer,production_designer	tt1692316,tt1297818,tt2106692
Tom Hankins			visual_effects,animation_department,visual_effects	tt1382720,tt3370776,tt2404307
Tom Hankins			director,writer,director	tt2492808

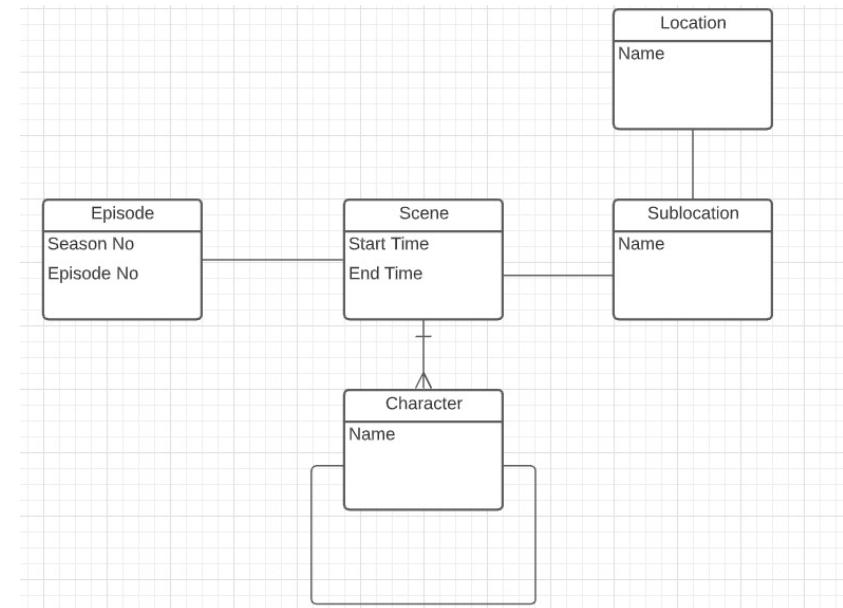
That was cool!

# *Game of Thrones*

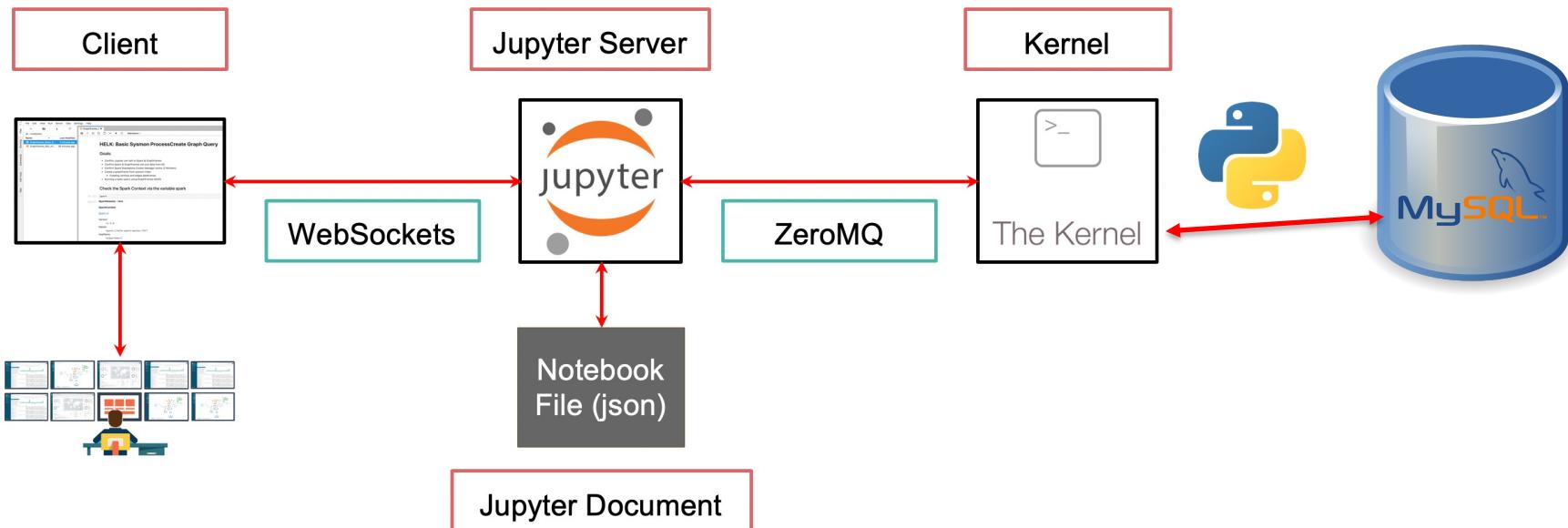
## *(Non-Programming Track Example)*

# Games of Thrones

- There are some interesting web sites and projects that analyze GoT.  
I used <https://jeffreylancaster.github.io/game-of-thrones/>
- I had to download, process and cleanup, and load into a DB.
- Structure
  - Seasons contain episodes.
  - Episodes contain scenes.
  - Scenes happen in a location.
  - Characters appear in scenes.
  - Character have interactions/relations
    - Marry
    - Kill
    - Protect
    - ... ...



# Data Analysis and Visual with Python

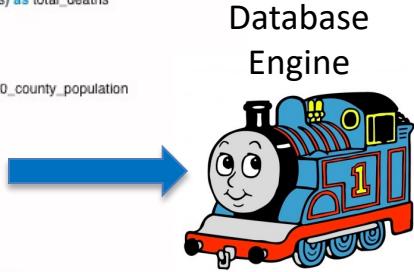


- Jupyter notebook contains:
  - Text in Markdown cells to explain the document, content, conclusions, ... ...
  - Python code and library calls to access application/data, process information, etc.
- Code executes on in a “server process” and accesses a database.

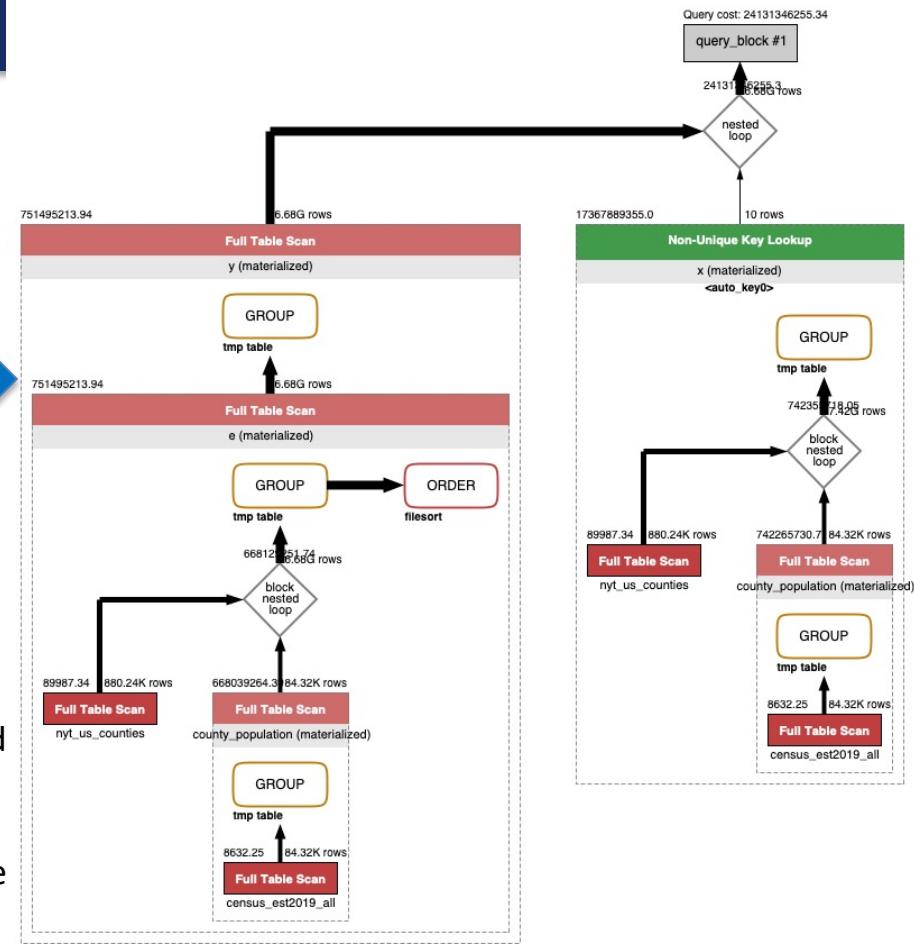
# Switch to Notebook

# Data Sciences Uses Complex Queries

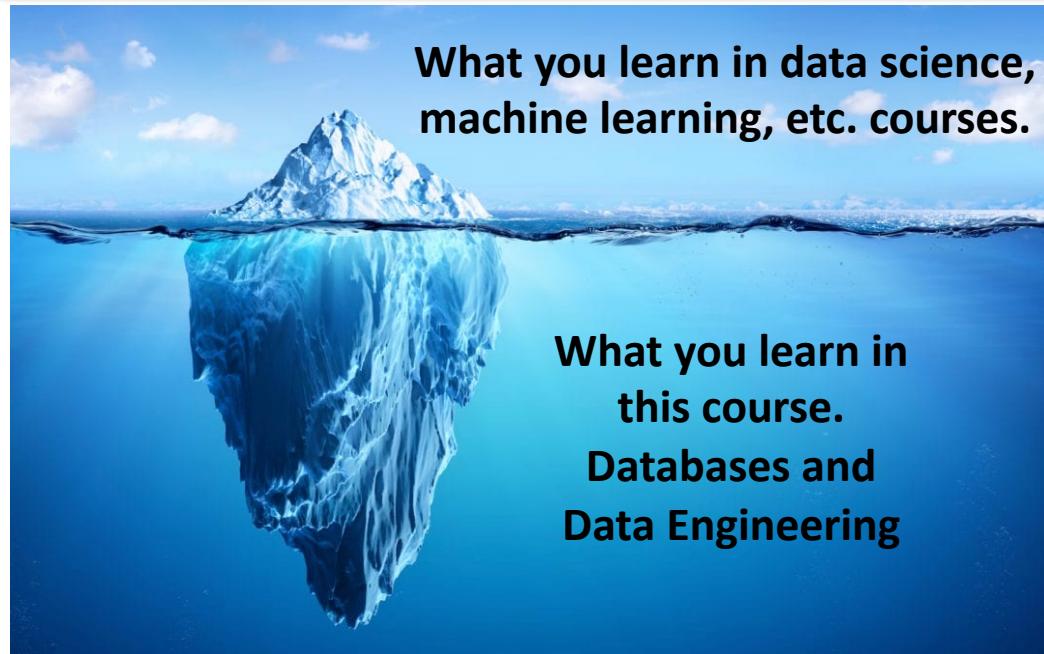
```
select c_year as 'Year', c_month as 'Month', state, cast(x.total_deaths as unsigned) as nyc_cumulative_deaths,  
cast(y.total_deaths as unsigned) as non_nyc_cumulative_deaths  
from  
(select  
year(date) as c_year, month(date) as c_month, county, state, max(deaths) as total_deaths  
from  
(select * from  
(SELECT  
concat(state, county) as fips, sum(popestimate2019) as estimated_2020_county_population  
FROM aaaaS21Examples.census_est2019_all  
group by fips) as county_population  
right join  
nyt_us_counties  
using(fips) as sub_result  
where county='New York City'  
group by c_year, c_month, county, state) as x  
join  
(select c_year, c_month, non_nyc, state, sum(total_deaths) as total_deaths from  
(select  
year(date) as c_year, month(date) as c_month, 'Non NYC' as non_nyc, state, max(deaths) as total_deaths  
from  
(select * from  
(SELECT  
concat(state, county) as fips, sum(popestimate2019) as estimated_2020_county_population  
FROM aaaaS21Examples.census_est2019_all  
group by fips) as county_population  
right join  
nyt_us_counties  
using(fips) as sub_result  
where county !='New York City' and state='New York'  
group by c_year, c_month, county, state  
order by county, c_year, c_month) as e  
group by  
c_year, c_month, state) as y  
using(c_year, c_month, state);
```



- SQL seems complex. Once you learn SQL, it is more productive than hand writing code.
- We will also see that the engine performs sophisticated performance optimizations.



# Data Science (IEOR, Finance, ...) Require Data Engineering



**What you learn in data science,  
machine learning, etc. courses.**

**What you learn in  
this course.  
Databases and  
Data Engineering**

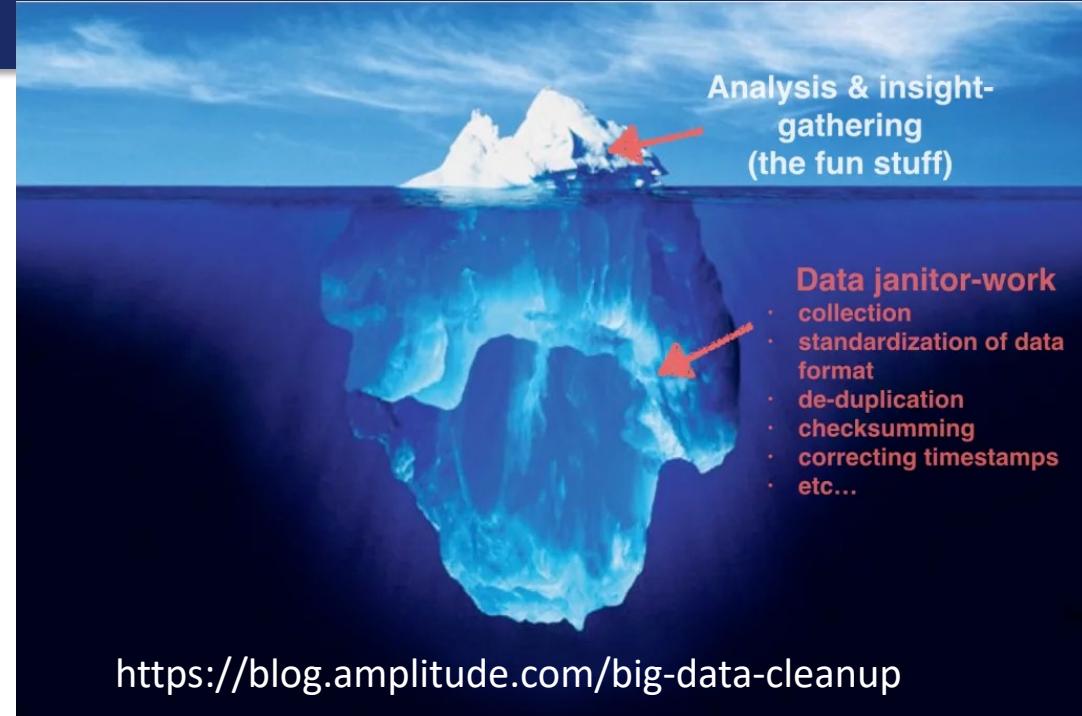
- 10% of the work is the cool data science, ML, ... ...
- 90% of the work is data engineering and query.

- You got a sense of the essence of database and data.
  - Query languages.
  - Schema.
  - Operations on data.
  - Scale and performance.

And ...

- Both programming and non-programming tracks will do work like what I just did.
- The non-programming track does a lot more of it.

# Data Cleansing



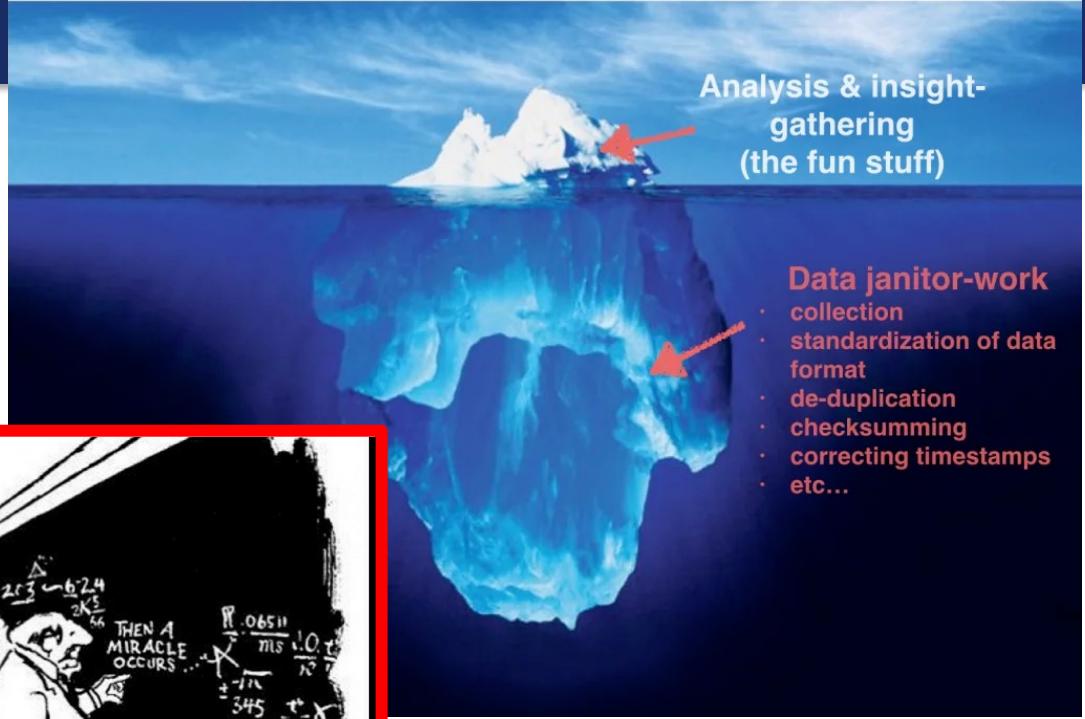
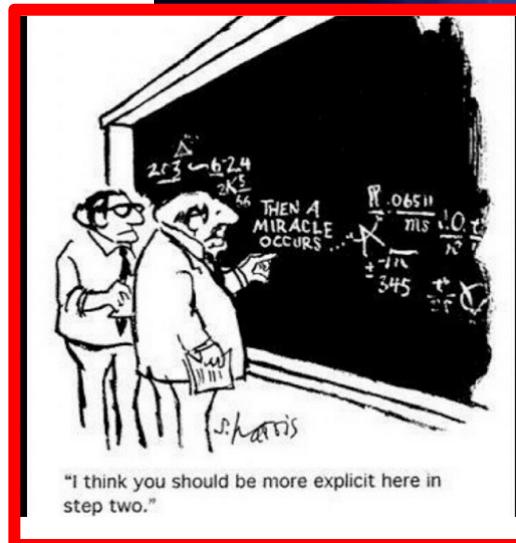
Database and data science classes **love to teach the fun stuff**  
queries, data modeling, machine learning, spooky math, algorithms, ... ...

# Data Cleansing

## Syllabus Topics

- Relational Foundations
- Overview (1 lecture)
- ER Model (2 lectures)
- Relational Model (4 lectures)
- Relational Algebra (2 lectures)
- SQL (5 lectures)
- Application Programming and Database APIs (1 lecture)
- Security (2 lectures)
- Normalization (2 lectures)
- Overview of Storage and Indexes (1 lecture)
- Overview of Query Optimization (1 lecture)
- Overview of Transaction Processing (1 lecture)
- Beyond Relational Foundations
- NoSQL (1 lecture)
- Data Preparation and Cleaning (1 lecture)
- Graphs (1 lecture)
- Object-Relational Databases (2 lectures)
- Cloud Databases (1 lecture)

## Recommended Syllabus



I place more emphasis on data cleansing and refactoring than other sections of W4111.

# 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.
  - I can cover in OHs, recitation, Q&A, Piazza, ... ...
- 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.

## Core Concepts:

- Data
- Databases
- Database management systems.
- Applications, application architectures.
- Roles.
- History of databases.
- etc.

Covering these topics is not a good use of lecture time:

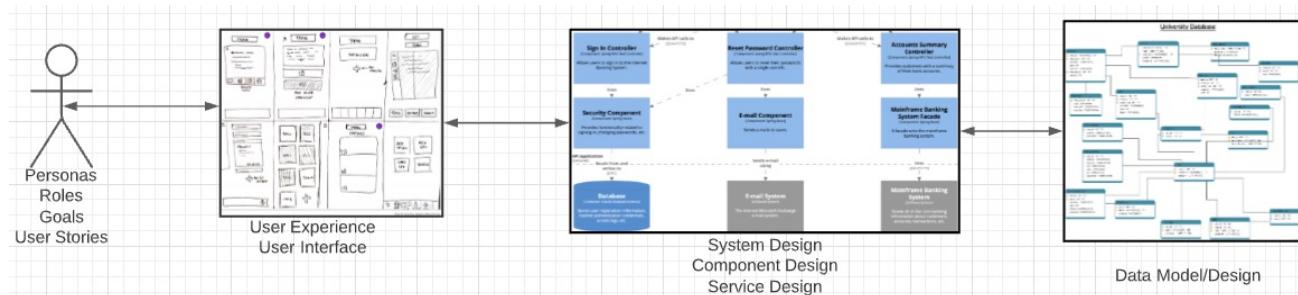
- Read the slides for chapter 1 from the recommended textbook.  
(<https://www.db-book.com/db7/slides-dir/PPTX-dir/ch1.pptx>)
- We will cover the material through questions on the reading, recitation, Piazza questions, ... ...

## Database design, Entity-Relationship Model (Part 1)

- Entity, relationship, attribute.
- Database design process.
- ER-Model and diagrams.

# Problem Statement

- We must build a system that supports academic operations in a university.
- There are two major domains:
  - Interactive operations, e.g. register, choose class, assign grade, ... ...
  - Insight and reporting, e.g. enrollment trends, overloaded resources, ... ...
- We will design, develop, test and deploy the system iteratively and continuously.
- There are four core domains.



- The processes are iterative, with continuous extension and details.

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



# 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 see slides with this formatting, they 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.



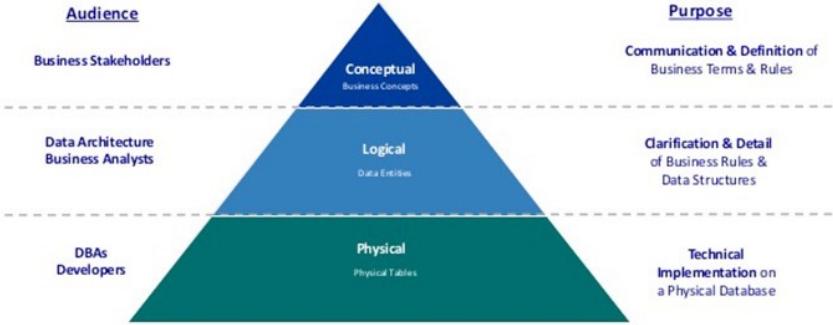
## Design Phases (Cont.)

- 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://ehkioya.com/conceptual-logical-physical-database-modeling/>

## Levels of Data Modeling

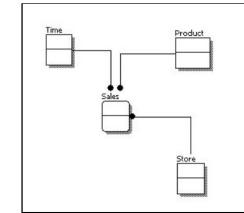


- It is easy to get carried away with modeling. You can spend all your time modeling and not actually build the schema.
- We will use the approaches in class.
- Mostly to understand concepts and patterns.

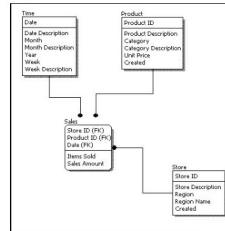
<https://www.1keydata.com/datawarehousing/data-modeling-levels.html>

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

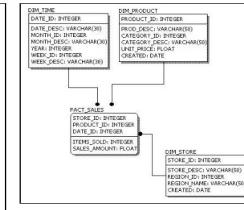
Conceptual Model Design



Logical Model Design



Physical Model Design



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



# ER model -- Database Modeling

- The ER data mode 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 2021

- 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:  
 $\text{instructor} = (\text{ID}, \text{name}, \text{salary})$   
 $\text{course} = (\text{course\_id}, \text{title}, \text{credits})$
- 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*

76766	Crick
45565	Katz
10101	Srinivasan
98345	Kim
76543	Singh
22222	Einstein

*instructor*

98988	Tanaka
12345	Shankar
00128	Zhang
76543	Brown
76653	Aoi
23121	Chavez
44553	Peltier

*student*



# Representing Entity sets in ER Diagram

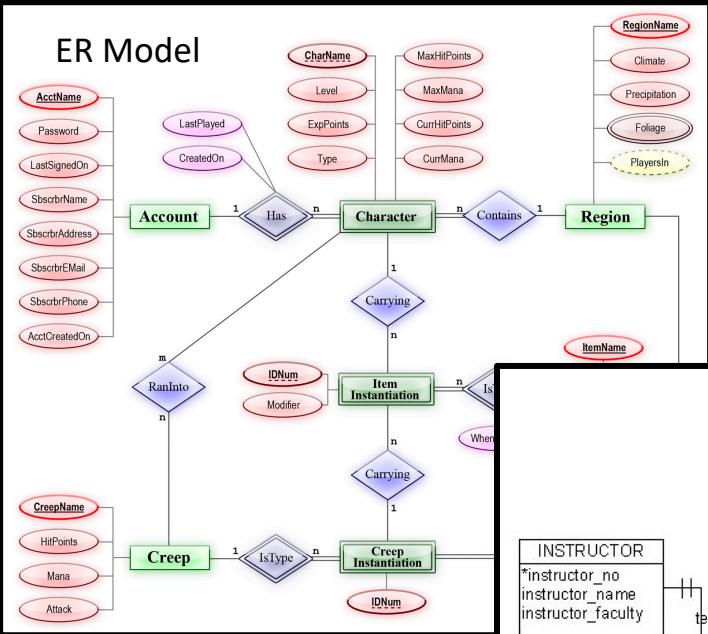
- Entity sets can be represented graphically as follows:
  - Rectangles represent entity sets.
  - Attributes listed inside entity rectangle
  - Underline indicates primary key attributes

<i>instructor</i>
<u>ID</u>
<i>name</i>
<i>salary</i>

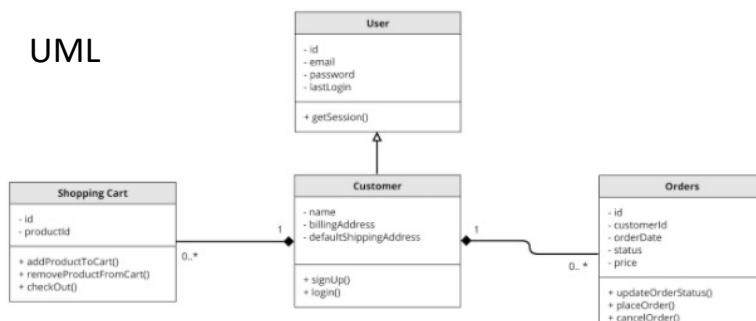
<i>student</i>
<u>ID</u>
<i>name</i>
<i>tot_cred</i>

# Visual Notation – Many Notations

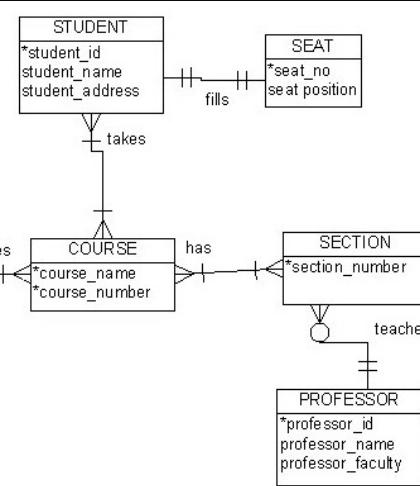
ER Model



UML



Crow's Foot



- “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.

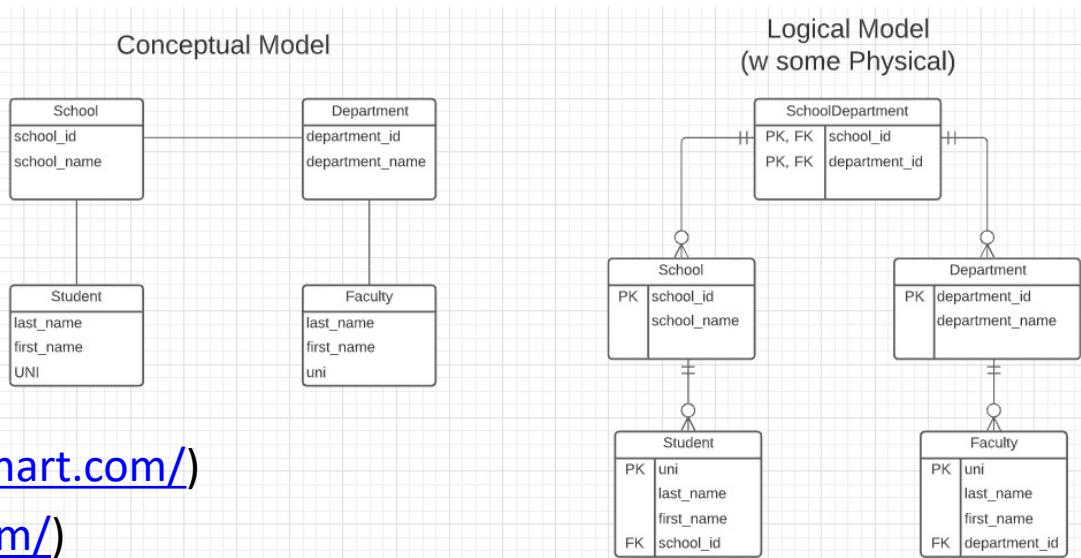
# First Iteration/Step for University Data Model

- 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/>)

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





# Entity Sets -- *instructor* and *student*

76766	Crick
45565	Katz
10101	Srinivasan
98345	Kim
76543	Singh
22222	Einstein

*instructor*

98988	Tanaka
12345	Shankar
00128	Zhang
76543	Brown
76653	Aoi
23121	Chavez
44553	Peltier

*student*

DFF Comments: Just a reminder of a previous slide.



# Relationship Sets

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

Example:

44553 (Peltier)                  advisor                  22222 (Einstein)  
*student entity*                  relationship set                  *instructor 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:

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

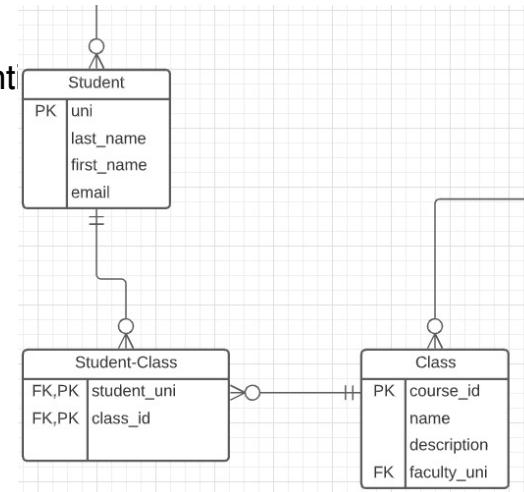
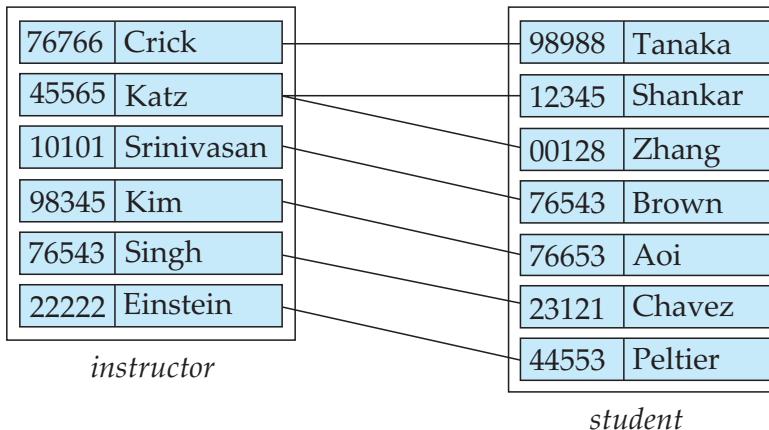
## DFF Comments:

- 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 Comments:

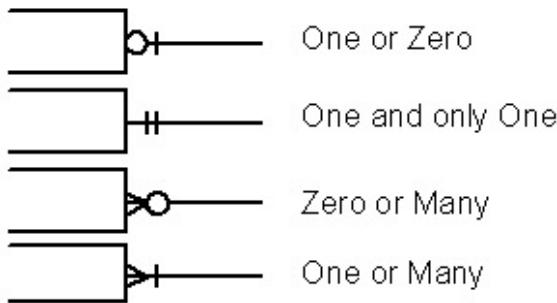
- Nobody draws the diagrams this way, but ...
- Sometimes thinking this way helps understand other ways to depict the concept.

(dff9, W4111)

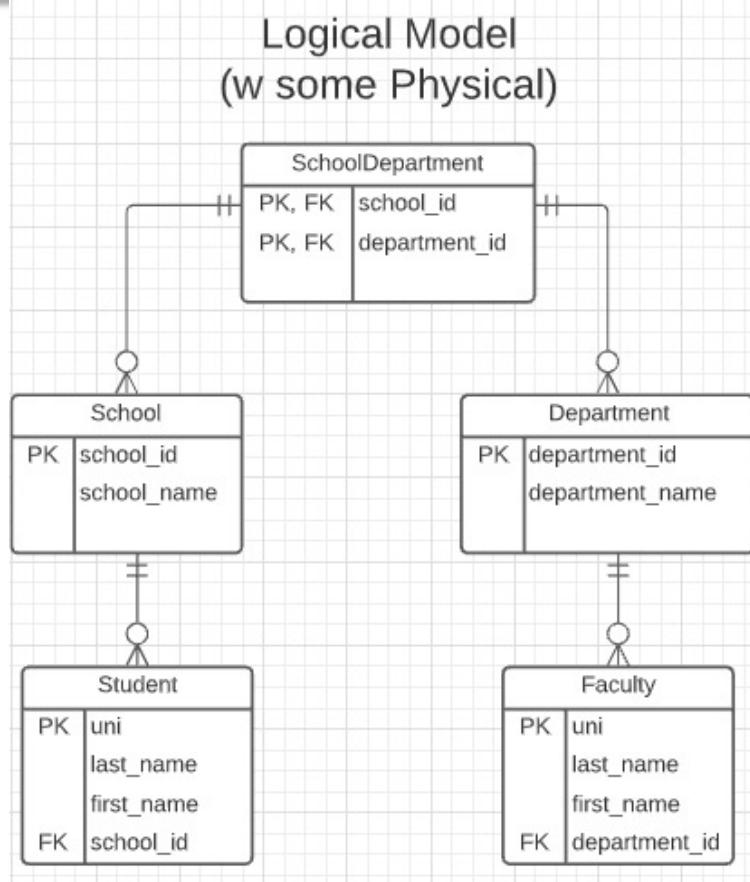
# Notation has Precise Meaning

- 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.



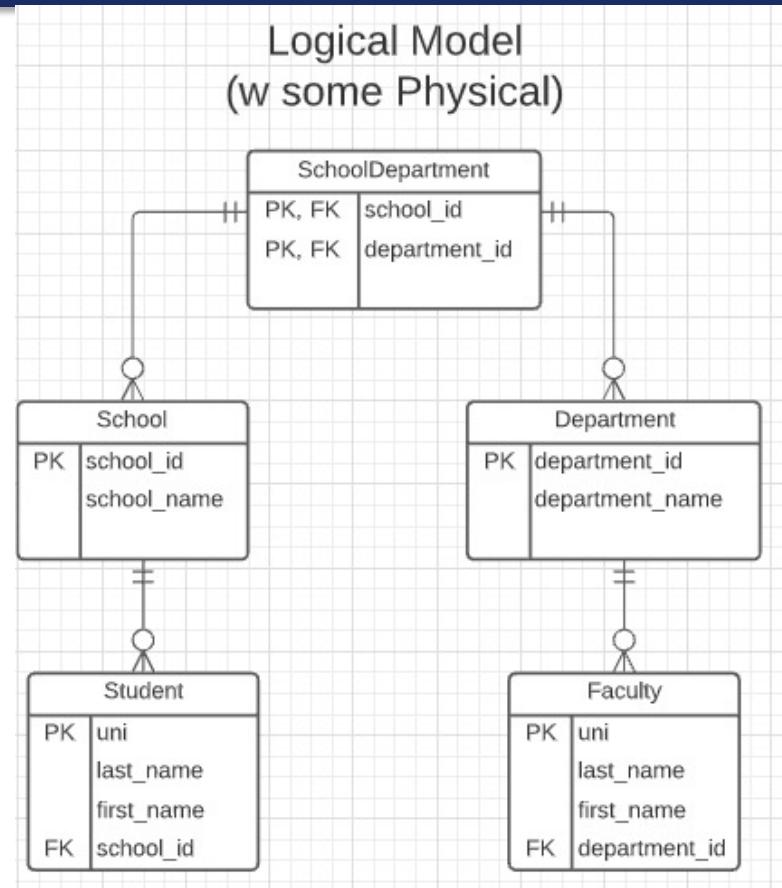
# What Does this Mean? Let's Get Started

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

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

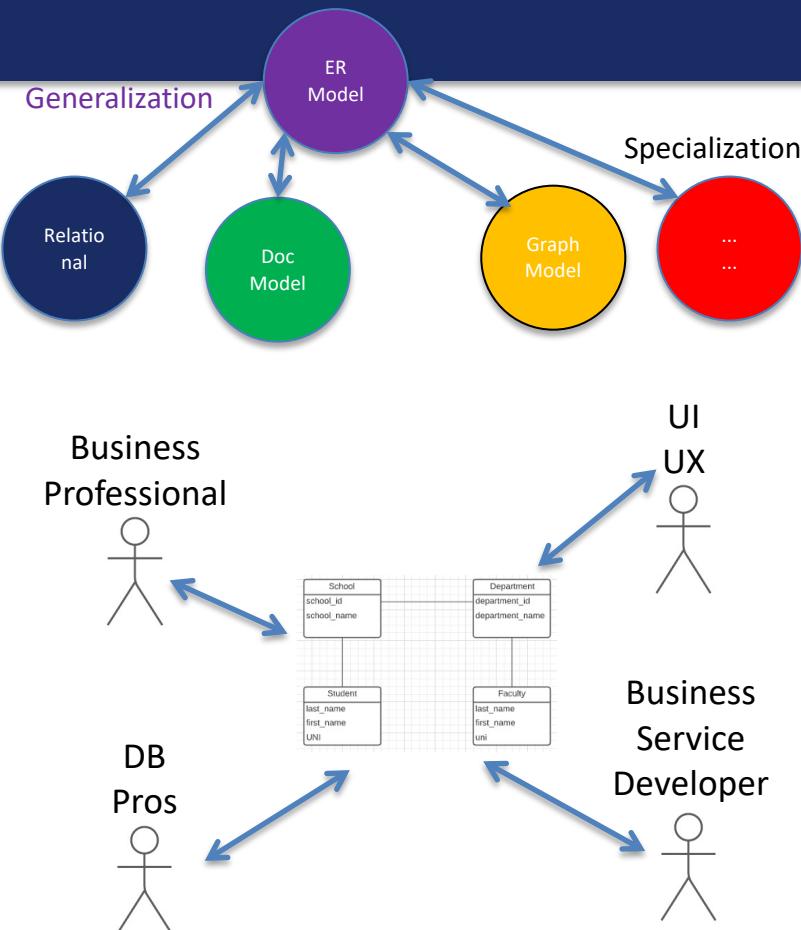
- The ER data mode was developed to facilitate database design by allowing specification of an **enterprise schema** that represents the overall logical structure of a database.
- The ER data model employs three basic concepts:
  - entity sets,
  - relationship sets,
  - attributes.
- The ER model also 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://pctechnicalpro.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.

# The Relational Model



# Relational Model

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



Ted Codd  
Turing Award 1981

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

- 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

<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

Diagram illustrating the components of the table:

- Attributes (or columns):** Four arrows point from the text "attributes (or columns)" to the column headers *ID*, *name*, *dept\_name*, and *salary*.
- Tuples (or rows):** Two arrows point from the text "tuples (or rows)" to the first two rows of the table data.



# 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

## 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

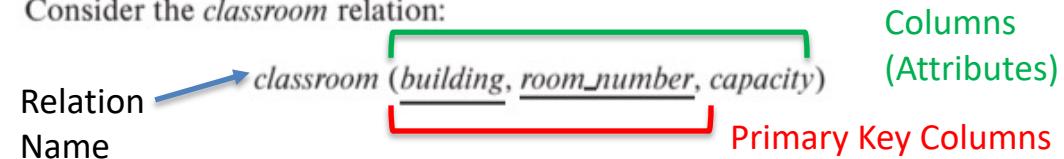
## Classroom relation

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:



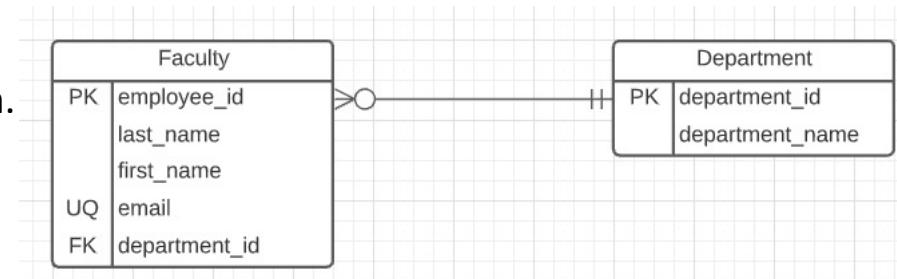
- The primary key is a *composite key*. Neither column is a key (unique) by itself.
- Keys are statements about all possible, valid tuples and not just the ones in the relation.
  - Capacity is unique in this specific data, but clearly not unique for all possible data.
  - In this domain, there cannot 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.

# Start Building our Datamodel

- First two relations:
  - Departments(title, id)
  - Faculty(employee\_id, first\_name, last\_name, email, **department\_id**)
- Note:
  - *email* is almost certainly a *candidate key*.
  - A primary key is chosen arbitrarily. Why did I use *employee\_id*?
    - Heuristic: People may change emails. I may lose “identity” in the data over time.  
Same person but two different emails.
    - *employee\_id* is *immutable*. I cannot change.
- ER Diagram:
  - Note: “UQ” is my non-standard convention.
  - Relationship:
    - A faculty “has” exactly one department.
    - A department has 0, 1 or many faculty.



# Start Building our Datamodel

OK, that is interesting, but ... ...

1. Where does the data come from?
  1. Users enter and update data through an application.
  2. You can *import* pre-existing data.
  3. You can generate test data.
2. How do you “manipulate” the data (CRUD)
  1. Create.
  2. Retrieve.
  3. Update.
  4. Delete.

# Import the Data – I screen-scraped the Website

The screenshot shows a web browser with the 'Elements' developer tool open, displaying the HTML structure of a department listing page. The page itself shows a grid of department cards, each with a title, a brief description, and a link. The developer tools show the detailed HTML for one of these cards, including its ID, class, and various attributes like 'summary' and 'path'.

```
-display paragraph--view-mode--default anchored">
  <div class="views-conditional-container container">
    ::before
    <div class="field field--name-field-cu-view-content field--type-viewfield field--label-hidden">
      <div class="viewfield-wrapper">
        <div class="viewfield">
          <div id="services-listing" class="view view-cu-services view-id-cu_services view-display-id-block_1 js-view-dom-id-51eb873b2f79adff8a8c4436adc2f4cd5197f0d57ca8da6cf2c7ce3f7 db10d0be">
            <div class="view-content">
              <script>
                var services_data = [{"title": "Accounting Division", "summary": false, "path": "/content/accounting-division", "id": "33", "groups_id": [], "categories_id": ["7"], "departments_id": ["32"]}, {"url": "https://www.gsb.columbia.edu/faculty-research/divisions/accounting"}, {"title": "Anesthesiology Department", "summary": false, "path": "/content/anesthesiology-department", "id": "35", "groups_id": [], "categories_id": ["7"], "departments_id": ["33"]}, {"url": "http://www.cumc.columbia.edu/dept/anesthesiology"}, {"title": "Anthropology Department", "summary": false, "path": "/content/anthropology-department", "id": "102", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://anthropology.columbia.edu"}, {"title": "Applied Physics and Applied Mathematics Department", "summary": false, "path": "/content/applied-physics-and-applied-mathematics-department", "id": "103", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://apam.columbia.edu"}, {"title": "Architecture (Barnard College) Department", "summary": false, "path": "/content/architecture-barnard-college", "id": "108", "groups_id": [], "categories_id": ["7", "9"], "departments_id": []}, {"url": "https://architecture.barnard.edu"}, {"title": "Art History and Archaeology Department", "summary": false, "path": "/content/art-history-and-archaeology-department", "id": "100", "groups_id": [], "categories_id": ["7"]}, {"url": "http://www.columbia.edu/cu/arthistory"}, {"title": "Astronomy and Astrophysics Department", "summary": false, "path": "/content/astronomy-and-astrophysics-department", "id": "116", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://www.astro.columbia.edu"}, {"title": "Biological Sciences Department", "summary": false, "path": "/content/biological-sciences-department", "id": "162", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://biology.columbia.edu"}, {"title": "Biomedical Engineering Department", "summary": false, "path": "/content/biomedical-engineering-department", "id": "163", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://bme.columbia.edu"}, {"title": "Biomedical Informatics, Department", "summary": false, "path": "/content/biomedical-informatics-department", "id": "165", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://bmi.columbia.edu"}, {"title": "Biostatistics Department", "summary": false, "path": "/content/biostatistics-department", "id": "167", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "https://www.dmbi.columbia.edu"}, {"title": "Chemical Engineering Department", "summary": false, "path": "/content/chemical-engineering-department", "id": "267", "groups_id": [], "categories_id": ["7"], "departments_id": []}, {"url": "http://cheme.columbia.edu"}, {"title": "Civil Engineering and Engineering Mechanics Department", "summary": false, "path": "/content/civil-engineering-and-engineering-mechanics-department", "id": "282", "groups_id": [], "categories_id": ["47"]}, {"url": "http://civil.columbia.edu"}, {"title": "Classics Department", "summary": false, "path": "/content/classics"}, {"url": "http://classics.columbia.edu"}]
```

# Import the Data – I screen-scraped the Website

Switch to the notebook.  
(cu\_info)

```
-display paragraph--view-mode--default anchored">
  <div class="views-conditional-container container">
    ::before
    <div class="field field--name-field-cu-view-content field--type-viewfield field--label-hidden">
      <div class="viewfield-wrapper">
        <div class="viewfield">
          <div id="services-listing" class="view view-cu-services view-id-cu_services view-display-id-block_1 js-view-dom-id-51eb873b2f79adff8a8c4436adc2f4cd5197f0d57ca8da6cf2c7ce3f7">
            <div class="view-content">
              <script>
                var services_data = [{"title": "Accounting Division", "summary": false, "path": "\content\accounting-division", "id": "33", "groups_id": [], "categories_id": [7], "departments_id": []}, {"title": "Anesthesiology Department", "summary": false, "path": "\content\anesthesiology-department", "id": "13", "groups_id": [1], "categories_id": [1], "departments_id": []}, {"title": "Anthropology (Barnard College)", "summary": false, "path": "\content\anthropology-barnard-college", "id": "167", "groups_id": [1], "categories_id": [7], "departments_id": [43], "url": "https://www.mailman.columbia.edu/become-student/departments/biostatistics"}, {"title": "Anthropology Department", "summary": false, "path": "\content\anthropology-department", "id": "167", "groups_id": [1], "categories_id": [7], "departments_id": [43], "url": "https://www.mailman.columbia.edu/become-student/departments/biostatistics"}, {"title": "Applied Physics and Applied Mathematics", "summary": false, "path": "\content\applied-physics-and-applied-mathematics", "id": "267", "groups_id": [1], "categories_id": [15], "departments_id": [15], "url": "http://chem.columbia.edu"}, {"title": "Architecture (Barnard College)", "summary": false, "path": "\content\architecture-barnard-college", "id": "282", "groups_id": [15], "categories_id": [47], "departments_id": [15], "url": "http://arch.columbia.edu"}, {"title": "Art History (Barnard College)", "summary": false, "path": "\content\art-history-barnard-college", "id": "282", "groups_id": [15], "categories_id": [47], "departments_id": [15], "url": "http://arch.columbia.edu"}, {"title": "Biostatistics Department", "summary": false, "path": "\content\biostatistics-department", "id": "167", "groups_id": [1], "categories_id": [7], "departments_id": [43], "url": "https://www.mailman.columbia.edu/become-student/departments/biostatistics"}, {"title": "Chemical Engineering Department", "summary": false, "path": "\content\chemical-engineering-department", "id": "267", "groups_id": [1], "categories_id": [7], "departments_id": [15], "url": "http://chem.columbia.edu"}, {"title": "Civil Engineering and Engineering Mechanics Department", "summary": false, "path": "\content\civil-engineering-and-engineering-mechanics-department", "id": "282", "groups_id": [15], "categories_id": [47], "departments_id": [15], "url": "http://arch.columbia.edu"}, {"title": "Classics Department", "summary": false, "path": "\content\classics-department", "id": "282", "groups_id": [15], "categories_id": [47], "departments_id": [15], "url": "http://arch.columbia.edu"}]
```

# Comments

- You probably have no idea what I just did.
- I wanted to give you a feel for data engineering. You saw some steps:
  - Find data from various sources.
  - Convert data into a format you can process.
  - Perform transformation, links, etc.
  - Save the data.This is often (historically) called Extract-Transform-Load (ETL)
- You are now ready to start integrating with other data.
- These steps enable analysis.
- These are some of things we will do in the non-programming track to enable more complex data science, etc.

# Generate Test Data – I Used to Tool

mockaroo SCHEMAS 6 DATASETS 5 SCENARIOS APIS PROJECTS 

My Schemas / university\_faculty

## university\_faculty

Save Changes

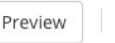
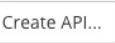
Field Name	Type	Options
employee_id	GUID	blank: 0 %  
first_name	First Name	blank: 0 %  
last_name	Last Name	blank: 0 %  
email	Email Address	blank: 0 %  
department_id	Dataset Column	university_departments  department_id  random  blank: 0 %  

Add another field

---

# Rows: 1000 Format: CSV  Line Ending: Unix (LF)  Include:  header  BOM

---

  |  



# 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:  $\sigma$
  - project:  $\Pi$
  - union:  $\cup$
  - set difference:  $-$
  - Cartesian product:  $\times$
  - rename:  $\rho$



# 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.

# 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:
  - Does have the sample data from the textbook.
  - You can also upload new data. (Show how and example)
- Some queries:
  - $\sigma \text{ title='Computer Science Department'} \vee \text{title='Accounting Division'} \text{ (departments)}$
  - $\pi \text{ last\_name, email } \text{ (faculty)}$
  - $\pi \text{ last\_name, email } ($   
 $\sigma \text{ department\_id=165} \vee \text{department\_id=1010 } \text{ (faculty)}$   
)

# SQL



# 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.

# SQL Language Statements

The core SQL language statements are:

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

$$\begin{aligned} \pi_{\text{last\_name}, \text{email}} ( \\ & \sigma_{\text{department\_id}=165 \vee} \\ & \text{department\_id}=1010 (\text{faculty}) \\ ) \\ = \\ \text{SELECT last\_name, email FROM faculty} \\ \text{WHERE department\_id=165 or} \\ \text{department\_id=1019} \end{aligned}$$



# 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*

$(A_1 D_1, A_2 D_2, \dots, A_n D_n,$   
 $\text{(integrity-constraint}_1\text{)},$   
 $\dots,$   
 $\text{(integrity-constraint}_k\text{)})$

- 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))
```

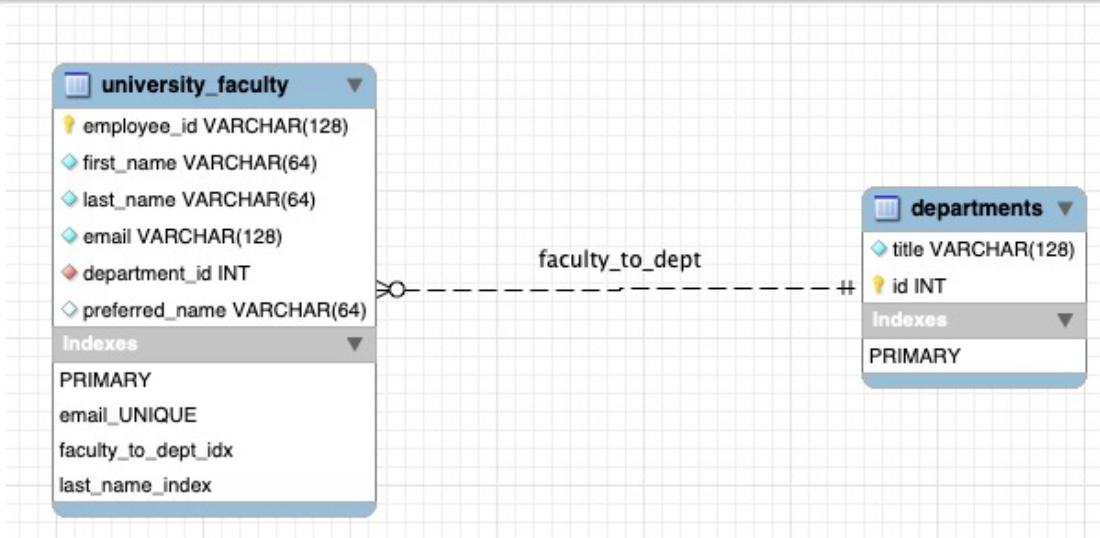
# SQL Create Table for our Datamodel

```
CREATE TABLE `departments` (
  `title` varchar(128) NOT NULL,
  `id` int NOT NULL,
  PRIMARY KEY (`id`)
) ENGINE=InnoDB
  DEFAULT CHARSET=utf8mb4
  COLLATE=utf8mb4_0900_ai_ci;
```

```
CREATE TABLE `university_faculty` (
  `employee_id` varchar(128) NOT NULL,
  `first_name` varchar(64) NOT NULL,
  `last_name` varchar(64) NOT NULL,
  `email` varchar(128) NOT NULL,
  `department_id` int NOT NULL,
  `preferred_name` varchar(64) DEFAULT NULL,
  PRIMARY KEY (`employee_id`),
  UNIQUE KEY `email_UNIQUE` (`email`),
  KEY `faculty_to_dept_idx` (`department_id`),
  KEY `last_name_index` (`last_name`),
  CONSTRAINT `faculty_to_dept` FOREIGN KEY (`department_id`) REFERENCES `departments` (`id`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci;
```

- The ` (back tick) is optional in most cases.
- You can look up the various column types.
  - INT is an integer.
  - VARCHAR(N) is a variable length string with maximum length 128.
- DEFAULT CHARSET, ENGINE, COLLATE are MySQL specific and you can usually ignore.
- NOT NULL is a Constraint. Constraints are core parts of SQL and we will cover in detail.
- Keys/Constraints/Index
  - You can see PRIMARY KEY definition.
  - UNIQUE KEY is similar to Candidate Key. UNIQUE KEY and the columns NOT NULL is a Candidate Key.
  - Foreign key defines source columns and reference columns.
  - “Key” sometimes means “Index” and is just used for performance.
- Some databases conflate the concepts of Key and Index. We will cover both keys and indexes in detail.

# Crow's Foot Like Physical Diagram



- It is possible to draw a physical model in an ER Diagram tool and generate the CREATE TABLE statements.
  - I never do.
  - I find it more trouble than it is worth.
- I will *reverse engineer* an existing database to understand its schema.

# 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.