

# CMSC424: Database Design

## Introduction/Overview

Professor: Pete Keleher  
keleher@umd.edu

---

### Today

- ▶ Administrivia
- ▶ Motivation: Why study databases ? What are databases ?
- ▶ Current Industry Outlook
- ▶ A typical DBMS at a glance

# Logistics

- ▶ Professor: Peter Keleher
  - 5146 Iribe Bldg
  - [keleher@umd.edu](mailto:keleher@umd.edu)
  - Class Webpage:
  - <http://ceres.cs.umd.edu/424>
- ▶ Communication:
  - Piazza
  - Office hours
  - Email to me: **include 424 in subject** as a last resort.
  - **Do not** message me on ELMS, I will not read it.

# Logistics

## Grading

All grades will be on [grades.cs.umd.edu](http://grades.cs.umd.edu).

### 28% Programming Assignments

We have 7 graded programming assignments. Each is worth 4% of the grade (1a + 1b are each 2%). All are due **Sunday at midnight**.

### 60% Exams

We have (3) exams:

- Exam 1 is 20%.
- Exam 2 is 20%.
- Exam 3 is 20%.
- *there is no final exam*

### 12% Weekly Homeworks

12 weekly homeworks:

- Each is worth 1%.
- All are due **Monday at midnight**.

# Logistics

- ▶ Grading
  - Whole class is curved: avg is B min, stdev up or down for A, C
  - Approximate cut-offs *last year* (not guaranteed)
    - 85+: A-
    - 75+: B-
    - 65+: C-
    - 60-: D/F
- ▶ Most had 40+ points (out of 50) on non-exams last year
  - *Must average a passing grade on the total exam score*

# Logistics

- ▶ Web site: <https://ceres.cs.umd.edu/424>
- ▶ Discussion: <https://piazza.com/class/m0390wa7rku3hn>
- ▶ Grades: <https://grades.cs.umd.edu>
- ▶ Gradescope: <https://www.gradescope.com/courses/424744>
  - homeworks, assignment submissions, graded exams
- ▶ Office Hours
  - Pete (me) IRB 5146, Tues 1:30 - 3:30 for lectures, exams, logistics
  - TAs (hours TBD):
    - Shayan Shabihi (“shayan”)
    - Anastasios Toumazatos (“tasos”)
- ▶ ELMS
  - *Nope!*

## Some To-Dos

- ▶ Sign up for Piazza !
  - If not already added
- ▶ Set up the computing environment (Assign. 0), and make sure you can run Docker, PostgreSQL, etc.
- ▶ Upcoming:
  - Homework 1 (due Monday),
  - Assign 0: Environment. Setting up Docker and PostgreSQL
    - due *next* Sunday midnight, but not graded/no submission
  - Assign 1: SQL. (*third* Sunday, midnight)

## Motivation: Data Overload

- ▶ Explosion of data, in pretty much every domain
  - Sensing devices and sensor networks that can monitor everything 24/7 from temperature to pollution to vital signs
  - Increasingly sophisticated smart phones
  - Internet, social networks makes it easy to publish data
  - Scientific experiments and simulations produce astronomical volumes of data
  - Internet of Things
  - **Datafication**: taking all aspects of life and turning them into data (e.g., what you like/enjoy turned into a stream of your "likes")
- ▶ How to handle that data? How to extract interesting actionable insights and scientific knowledge?
- ▶ Data volumes expected to get much worse

## Four V's of Big Data

- ▶ Increasing data Volumes
  - Scientific data: 1.5GB/genome -- can be sequenced in .5 hrs; LHC generates 100TB of data a day
  - 500M tweets per day (as of 2013)
  - As of 2012: 2.5 Exabytes of data created every day
  - EBay: Two data warehouses with 7.5PB and 40PB
  - Walmart: 583 terabytes of sales and inventory data
  - FICO monitors 2.5 billion active accounts worldwide
- ▶ Variety:
  - Structured data, spreadsheets, photos, videos, natural text, ...

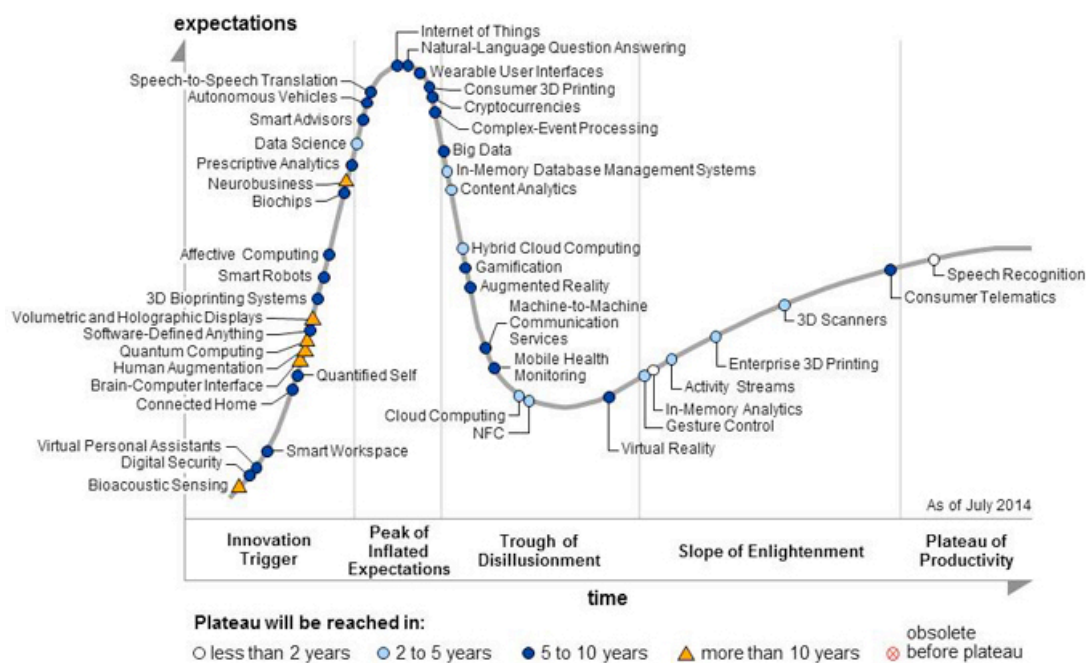
## Four V's of Big Data

- ▶ Velocity
  - Sensors everywhere -- can generate tremendous volumes of "data streams"
  - Real-time analytics requires data to be consumed as fast as it is generated
- ▶ Veracity
  - How do you decide what to trust? How to remove noise? How to fill in missing values?

# Big Data and Data Science to the Rescue

- ▶ Terms increasingly used synonymously: also data analytics, data mining, business intelligence
  - Loosely used for any process where interesting things are inferred from data
  - Google search: “How Big Data Will Change”
- ▶ Data scientist called the sexiest job of the 21st century
  - The term has becoming very muddled at this point

## Overhyped Words



## Is it all hype?

- ▶ No: Extracting insights and knowledge from data very important, and will continue to increase in importance
  - Big data techniques are revolutionizing things in many domains like Education, Food Supply, Disease Epidemics, ...
- ▶ But: it is not much different from what we, especially statisticians, have been doing for many years
- ▶ What is different?
  - Much more data is digitally available than was before
  - Inexpensive computing + Cloud + Easy-to-use programming frameworks = Much easier to analyze it
  - Often: large-scale data + simple algorithms > small data + complex algorithms
    - Changes how you do analysis dramatically

## Motivation: Data Overload

- ▶ How do we do anything with this data?
- ▶ Where and how do we store it ?
  - Disks are doubling every 18 months or so -- not enough
  - In many cases, the data is not actually recorded as it is; *summarized* first
- ▶ What if the disks crash ?
  - Very common, especially with 10,000's of disks
- ▶ How do we ensure "correctness" ?
  - What if the system crashes in the middle of an ATM transaction ?
    - Can't have money disappearing
  - What happens when a million people try to buy tickets to *<your favorite artist>'s concert* at the same time ?



## Motivation: Data Overload

- ▶ What to do with the data ? How to process/analyze it ?
  - text search ?
    - Very limited
  - “find the stores with the maximum increase in sales in last month”
    - We can’t expect the users to write Java programs
  - “how much time from here to Pittsburgh if I start at 2pm ?”
    - Data is there; more will be soon (GPS, live traffic data)
    - Requires predictive capabilities
  - Increasing need to convert “information” to “knowledge”: **Data mining**
    - “How should we replicate different movies?” (Netflix)
    - Find videos with this type of an event (say car break-ins)
    - Mine the “blogs” to detect “buzz”

## Motivation: Data Overload

- ▶ Speed !!
  - With TB’s of data, just finding something (even if you know what), is not easy
    - Reading a file with TB of data can take hours
  - Imagine a bank and millions of ATMs
    - How much time does it take you to do a withdrawal ?
    - The data is not local
- ▶ How do we guarantee the data will be there 10 years from now ?
- ▶ Privacy and security !!!
  - Every other day we see some database leaked on the web
    - identity fraud, influencing elections...
  - How to make sure different users’ data is protected from each other



## Why not use file systems ?

- ▶ Drawbacks of using file systems to store data:
  - Data redundancy and inconsistency
    - Multiple file formats, duplication of information in different files
  - Difficulty in accessing data
    - Need to write a new program to carry out each new task
  - Data isolation — multiple files and formats
  - Integrity problems
    - Integrity constraints (e.g., account balance  $> 0$ ) become “buried” in program code rather than being stated explicitly
    - Hard to add new constraints or change existing ones

## Why not use file systems ?

- ▶ Drawbacks of using file systems to store data:
  - Atomicity of updates
    - Failures may leave database in an inconsistent state with partial updates carried out
    - Example: Transfer of funds from one account to another should either complete or not happen at all
  - Concurrent access by multiple users
    - Concurrent access needed for performance
    - Uncontrolled concurrent accesses can lead to inconsistencies
      - Example: Two people reading a balance (say 100) and updating it by withdrawing money (say 50 each) at the same time
  - Security problems
    - Hard to provide user access to some, but not all, data

# Today

- ▶ Administrivia
  - Homework 1 due Monday 11:59 pm

## DBMSs to the Rescue

- ▶ Provide a systematic way to answer many of these questions...
- ▶ Aim is to allow easy management of high volumes of data
  - Storing , Updating, Querying, Analyzing ....
- ▶ What is a Database ?
  - A large, integrated collection of (mostly *structured*) data
  - Typically models and captures information about a real-world **enterprise**
    - **Entities** (*e.g. courses, students*)
    - **Relationships** (*e.g. John is taking CMSC 424*)
    - Usually also contains:
      - Knowledge of **constraints** on the data (*e.g. course capacities*)
      - **Business logic** (*e.g. pre-requisite rules*)
      - Encoded as part of the data model (preferable) or through external programs

# DBMSs to the Rescue

- ▶ Massively successful for *highly structured data*
  - Why ? Structure in the data (if any) can be exploited for ease of use and efficiency
    - If there is no structure in the data, hard to do much
    - Contrast managing emails vs managing photos
  - Much of the data we need to deal with is highly structured
  - Some data is *semi-structured*
    - E.g.: Resumes, Webpages, Blogs etc.
  - Some has complicated structure
    - E.g.: Social networks
  - Some has no structure
    - E.g.: Text data, Video/Image data etc.

## Structured vs Unstructured Data

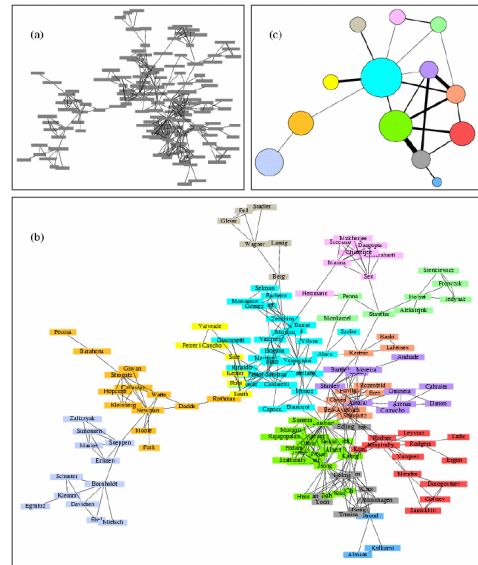
- Structured data often has very simple structures
  - E.g. Data that can be represented in tabular forms
- Significantly easier to deal with
- We will focus on such data for much of the class

Account		
bname	acct_no	balance
Downtown	A-101	500
Mianus	A-215	700
Perry	A-102	400
R.H	A-305	350

Customer		
cname	cstreet	ccity
Jones	Main	Harrison
Smith	North	Rye
Hayes	Main	Harrison
Curry	North	Rye
Lindsay	Park	Pittsfield

# Structured vs Unstructured Data

- ▶ Some data has a little **more complicated structure**
  - E.g graph structures
    - Map data, social networks data, the web link structure etc.
  - Can convert to tabular forms for storage, but may not be optimal
  - Queries often reason about graph structure
    - *Find my "Erdos number"*
    - *Suggest friends based on current friends*
  - Growing importance in recent years in a variety of domains: Biological, social networks, web...

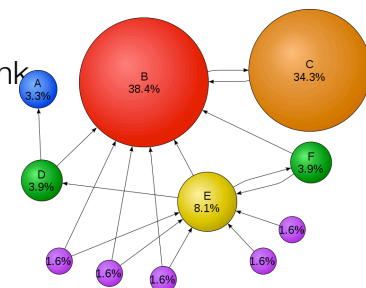


# Structured vs Unstructured Data

- ▶ Increasing amount of data in a **semi-structured format**
  - XML – Self-describing tags (HTML ?)
  - Complicates a lot of things
  - We will discuss this toward the end
- ▶ A huge amount of data is unfortunately **unstructured**
  - Books, WWW
  - Amenable to pretty much only text search... so far
    - Information Retrieval research deals with this topic
  - What about Google search ?
    - Google search is mainly successful because it uses link structure (in its original incarnation)
- ▶ Video ? Music ?
  - Can represent in DBMS's, but can't really operate on them

```

<Symbol>List</Symbol>
<Function>
<Symbol>List</Symbol>
<Symbol>Automatic</Symbol>
<Number>4.</Number>
</Function>
<Function>
<Symbol>List</Symbol>
<Symbol>Automatic</Symbol>
<Number>6.</Number>
</Function>
</Options>
</Notebook>
    
```



circle size == page importance == **pagerank**  
 more incoming links → higher pagerank  
 incoming links from important pages → higher pagerank

# DBMSs to the Rescue

- ▶ Massively successful for *highly structured data*
  - Two Key Concepts:
    - Data Modeling: Allows reasoning about the data at a high level
      - e.g. “emails” have “sender”, “receiver”, “...”
      - Once we can describe the data, we can start “querying” it
    - Data Abstraction/Independence:
      - Layer the system so that the users/applications are insulated from the low-level details

## DBMSs to the Rescue: Data Modeling

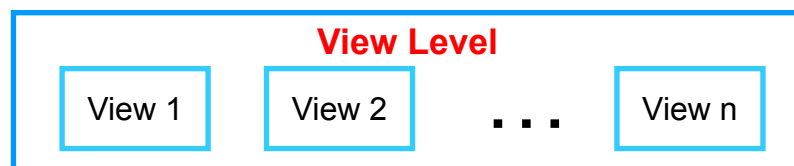
- ▶ Data modeling
  - **Data model**: A collection of concepts that describes how data is represented and accessed
  - **Schema**: A description of a specific collection of data, using a given data model
  - Some examples of data models that we will see
    - Relational, Entity-relationship model, XML, JSON...
    - Object-oriented, object-relational, semantic data model, RDF...
  - Why so many models ?
    - Tension between descriptive power and ease of use/efficiency
      - More powerful models → more data can be represented
      - More powerful models → harder to use, to query, and less efficient

# DBMSs to the Rescue: Data Abstraction

- ▶ Probably *the* most important purpose of a DBMS
- ▶ Goal: Hiding *low-level details* from the users of the system
  - Alternatively: the principle that
    - *applications and users should be insulated from how data is structured and stored*
  - Also called *data independence*
- ▶ Through use of *logical abstractions*

## Data Abstraction

What data users and application programs see ?



What data is stored ?

describe data properties such as data semantics, data relationships

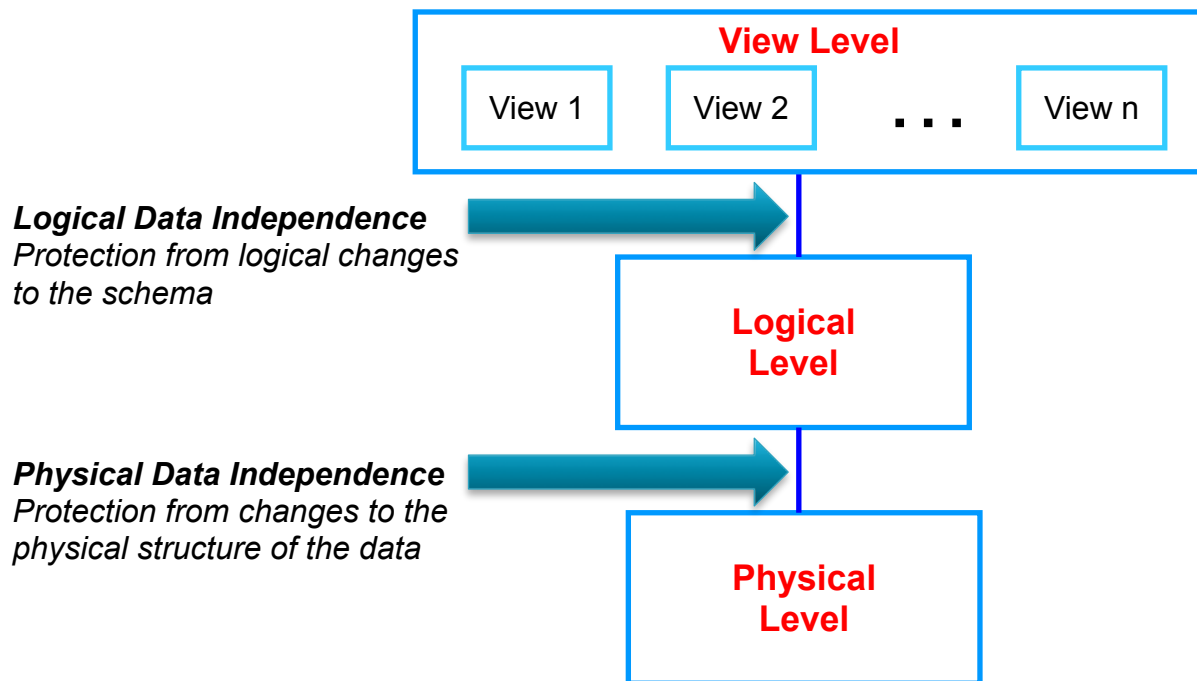


How data is actually stored ?

e.g. are we using disks ? Which file system ?



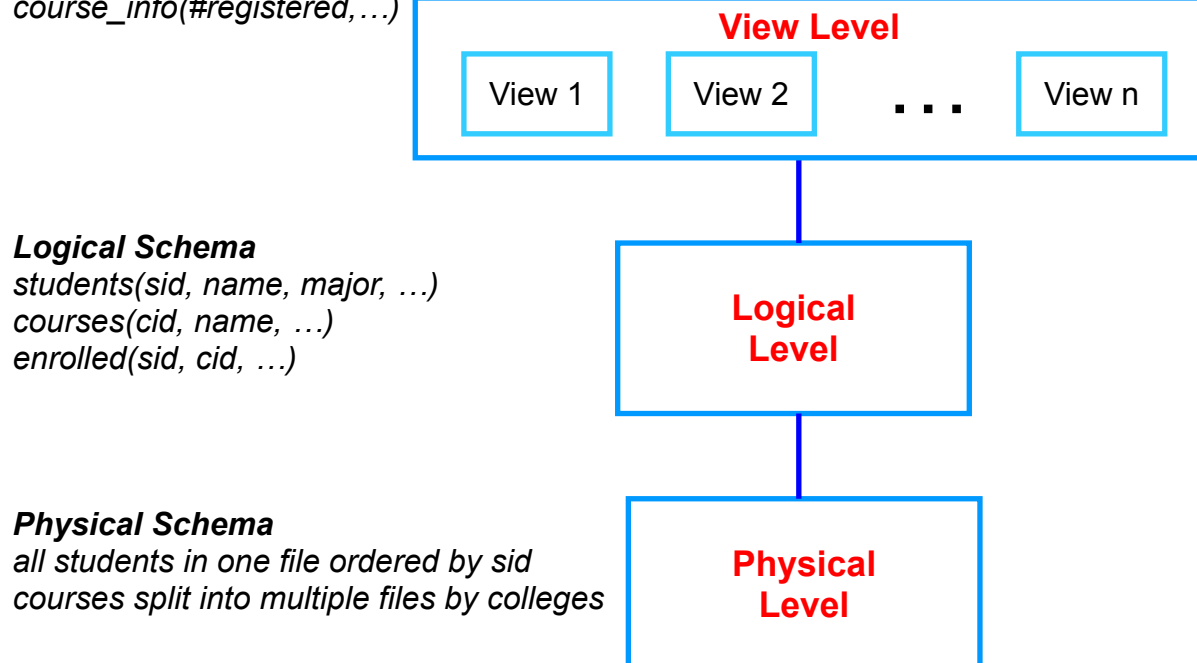
# Data Abstraction



## Data Abstractions: Example

### A "view" Schema

*course\_info(#registered, ...)*



# What about a Database System ?

- ▶ A DBMS is a software system designed to store, manage, facilitate access to databases
- ▶ Provides:
  - Data Definition Language (DDL)
    - For defining and modifying the schemas
  - Data Manipulation Language (DML)
    - For retrieving, modifying, analyzing the data itself
  - Guarantees about correctness in presence of failures and concurrency, data semantics etc.
- ▶ Common use patterns
  - Handling transactions (e.g. ATM Transactions, flight reservations)
  - Archival (storing historical data)
  - Analytics (e.g. identifying trends, **Data Mining**)

## Relational DBMS: SQL

- ▶ **SQL** (sequel): Structured Query Language
- ▶ **Data definition (DDL)**
  - **create table** *instructor* (  
                                  *ID*                  **char**(5),  
                                  *name*          **varchar**(20),  
                                  *dept\_name* **varchar**(20),  
                                  *salary*      **numeric**(8,2))
- ▶ **Data manipulation (DML)**
  - Example: Find the name of the instructor with ID 22222  
**select** *name*  
**from** *instructor*  
**where** *instructor.ID* = '22222'



# Current Industry Outlook

- ▶ Relational DBMSs
  - Oracle, IBM DB2, Microsoft SQL Server, Sybase
- ▶ Open source alternatives
  - MySQL, PostgreSQL, Apache Derby, BerkeleyDB (mainly a storage engine – no SQL), neo4j (graph data) ...
- ▶ Data Warehousing Solutions
  - Geared towards very large volumes of data and on analyzing them
  - Long list: Teradata, Oracle Exadata, Netezza (based on FPGAs), Aster Data (founded 2005), Vertica (column-based), Kickfire, Xtremedata (released 2009), Sybase IQ, Greenplum (eBay, Fox Networks use them)
  - Usually sell package/services and charge per TB of managed data
  - Many (especially recent ones) start with MySQL or PostgreSQL and make them parallel/faster etc..

# Web Scale Data Management, Analysis

- ▶ Ongoing debate/issue
  - Cloud computing seems to eschew DBMSs in favor of homegrown solutions
  - E.g. Google, Facebook, Amazon etc...
- ▶ MapReduce: A paradigm for large-scale data analysis
  - Hadoop: An open source implementation
  - Apache Spark: a better open source implementation
- ▶ Why ?
  - DBMSs can't scale to the needs, not fault-tolerant enough
    - These apps don't need things like transactions, that complicate DBMSs (???)
  - MapReduce favors Unix-style programming, doesn't require SQL
    - Try writing SVMs or decision trees in SQL
  - Cost
    - Companies like Teradata may charge \$100,000 per TB of data managed

# Current Industry Outlook

- ▶ Bigtable-like
  - Called “key-value stores”
  - Think highly distributed hash tables
  - Allow some transactional capabilities – still evolving area
  - PNUTS (Yahoo), **Apache Cassandra** (Facebook), Dynamo (Amazon), and many many others
- ▶ Mapreduce-like
  - Hadoop (open source), Pig (@Yahoo), Dryad (@Microsoft), Spark
  - Amazon EC2 Framework
  - Not really a database – but increasing declarative SQL-like capabilities are being added (e.g. HIVE at Facebook)
- ▶ Much ongoing research in industry and academia

# DBMS at a glance

- ▶ Data Models
  - Conceptual representation of the data
- ▶ Data Retrieval
  - How to ask questions of the database
  - How to answer those questions
- ▶ Data Storage
  - How/where to store data, how to access it
- ▶ Data Integrity
  - Manage crashes, concurrency
  - Manage semantic inconsistencies
- ▶ Not fully disjoint categorization !!

# Relational Query Languages

- ▶ Example schema:  $R(A, B)$
- ▶ Practical languages
  - SQL
    - select A from R where B = 5;
  - Datalog (sort of practical)
    - $q(A) :- R(A, 5)$
- ▶ Formal languages
  - Relational algebra  
 $\pi_A ( \sigma_{B=5} (R) )$
  - Tuple relational calculus  
 $\{ t : \{A\} \mid \exists s : \{A, B\} ( R(A, B) \wedge s.B = 5 ) \}$
  - Domain relational calculus
    - Similar to tuple relational calculus