

Logistics

Instructor:

- Prof. Bertram Ludaescher (ludaesch@ucdavis.edu)
– Office Hours: **T 8:30-9:30am, 3051 Kemper Hall**

Teaching Assistants:

- Meghan Raul (meghanraul@me.com)
- Harika Sabbella (hhsabbella@ucdavis.edu)
- Office Hours: TBD, 55 Kemper Hall

Discussion Sections:

- W 9-9:50am, 212 Wellman (TA / Instructor)
- **On demand** (not before **Week 2**)

ECS-165B

1

ECS-165B Database Systems

Home sites.google.com/site/ecs165bwinter2014/

Instructor:
* Prof. Bertram Ludaescher (ludaesch@ucdavis.edu)
Office Hours: T 8:30-9:30am, 3051 Kemper Hall

Teaching Assistants:
* Meghan Raul (meghanraul@me.com)
Office Hours: TBD, 55 Kemper Hall
* Harika Sabbella (hhsabbella@ucdavis.edu)
Office Hours: TBD, 55 Kemper Hall

MEETING	CNMF	TIME	ROOM	Staff
Lecture	01101	MWTF 3:10-4pm	118 Olson	B. Ludaescher
Discussion Section	01101	W 9-9:50am	212 Wellman	TA/Instructor

Class Mailing List:
Sign-up: TBD (Piazza or Google list)

Textbook
* *Fundamentals of Database Systems*, by Elmasri, Navathe, Addison-Wesley.
Online Textbook
* *Foundations of Databases*, Textbook by Abiteboul, Hull, Vianu.

Textbooks
* [GMUW09] *Database Systems: The Complete Book*, Garcia-Molina, Ullman, Widom, Prentice Hall; 2nd ed. (2009)
* [EN10] *Fundamentals of Database Systems*, Elmasri, Navathe, Addison-Wesley (2010)
* [SKS05] *Database System Concepts*, Silberschatz, Korth, Sudarshan (2005)

Online Textbook
* [JMV98] *Foundations of Databases*, Abiteboul, Hull, Vianu (1998)

Related Online Class (Stanford)
* *Introduction to Databases (Winter 2013)*

Database System:
We will mainly use PostgreSQL, available on the CSUF machines.
(We will also use other systems and paradigms: Log, MongoDB, MapReduce)

2

More Logistics

Class page:

- <https://sites.google.com/site/165bWinter2014>

Class Mailing List:

Sign-up:

- <https://piazza.com/ucdavis/winter2014/ecs165b/>

Textbooks

- [GMUW09] *Database Systems: The Complete Book*, Garcia-Molina, Ullman, Widom, Prentice Hall; 2nd ed. (2009)
- [EN10] *Fundamentals of Database Systems*, Elmasri, Navathe, Addison-Wesley (2010)
- [SKS05] *Database System Concepts*, Silberschatz, Korth, Sudarshan (2005)

ECS-165B

3

165A Course Topics (pre-req for 165B)

- Database Design, **E/R Model**
- Relational Model, **Relational Algebra**
- **SQL** (Structured Query Language)
- **Integrity Constraints**
- Storage structures, Indexing
- **Query Processing**
- Transactions
- Additional Topics & Current Trends
 - Logic/declarative queries ("**Datalog**")

ECS-165B

4

Prerequisite Screening

- ECS-165B prereq: ECS-165A
- Campus is now a bit stricter about screening prereqs
- If you haven't done 165A (or didn't do so well), you might still take 165B
- I will email individuals who failed the pre-req screening.

ECS-165B

5

165B Overview

- **Database Foundations** ("Theory")
 - Design Theory for Relational Databases (Normalization)
 - Recursive queries (Datalog, SQL-with-recursive)
 - Querying (XML) trees and graphs
- **Advanced & Hands-on Topics** ("Practice")
 - XML data management
 - Online Analytical Processing (OLAP), Data Cubes
 - Map-Reduce Parallelism Framework & "Big Data"
 - Other trends (NoSQL)

ECS-165B

6

165B Course Topics (tentative)

- **Database theory:**
 - recursive queries, data integration
 - DB design: normalization
 - data provenance
- **Semistructured data on the web (XML)**
 - DTDs, XML Schema
 - XPath
 - XQuery
 - XSLT
- **Advanced database topics**
 - OLAP (vs OLTP)
 - Big Data, parallel processing e.g. MapReduce
 - Specialized topics (Fusion Tables? NoSQL?)

ECS-165B

7

165B Course Topics

- Focus is on
 - **Foundations**
 - DB Theory (Normalization) and Datalog
 - **Practical experience** with SQL, XML,
 - We'll use PostgreSQL
 - A "real" (full-featured), scalable DBMS
 - Open source, available @CSIF and @home!
 - » Other systems: Oracle, SQL-Server, MySQL, SQLite, ...
 - » Embedded SQL (e.g. with Python)
 - ... and some other things ...
 - Map-Reduce @ Amazon
- **Individual Assignments (~3-4)**
- **Group Projects (~3)**

ECS-165B

8

Assignments (Examples)

- Individual Assignments:
 - Datalog, SQL With Recursive
 - XML DTDs, XPath, XQuery, XML Accelerator
 - Design Theory (Normalization)
 - Data Provenance
- Group Projects:
 - Graph Data Visualization Tool
 - XML to Relational Mapping
 - Big Data / Map-Reduce (Amazon EC2)

ECS-165B

9

ECS-199: Special Study for Advanced Undergraduates

- If you did very well in ECS-165A ...
- ... and you want to get some more additional hands-on experience!
- I can offer a couple (not too many) of these.
- Many possibilities:
 - Complex queries: on trees, graphs; Skylines; OLAP; temporal; spatial; ...
 - Harvesting data from the web (DBpedia, ...)
 - Analyzing (social) network data
 - Mobile DBs (Android)
 - ...

ECS-165B

10

<http://class2go.stanford.edu/db/Winter2013>

11

Grading and Policies

- **Grading:**
 - Base line:
 - 30% Individual Assignments
 - 30% Group Projects
 - 40% Exams
- **Academic Conduct**
 - Be polite
 - Don't cheat
- Ask when in doubt
- Make good use of the mailing-list/forum

ECS-165B

12

Why study databases / data management?

- Critical to business, government, science, culture, society, ...
- Determines success of many corporations (even their existence)
- Many tech companies built on data management (Google, Amazon, Yahoo!, Facebook, ...)
- ... or offer database products (Microsoft, IBM, Oracle)
- Database systems span major areas of computer science
 - Operating systems (file, memory, process management)
 - Theory (languages, algorithms, complexity)
 - Artificial Intelligence (knowledge-based systems, logic, search)
 - Software Engineering (application development)
 - Data structures (trees, hash-tables)
 - ... and the DB research community continues to be very active

ECS-165B

13

Lots of Data Everywhere

- From <http://en.wikipedia.org/wiki/Petabyte> :
- **History:** According to Kevin Kelly in *The New York Times*, "the entire [written] works of humankind, from the beginning of recorded history, in all languages" would amount to 50 petabytes of data.^[1]
- **Computer hardware:** *Teradata* Database 12 has a capacity of 50 petabytes of compressed data.^{[2][3]}
- **Telecoms:** AT&T has about 16 petabytes of data transferred through their networks each day.^[4]
- **Archives:** The *Internet Archive* contains about 3 petabytes of data, and is growing at the rate of about 100 terabytes per month as of March, 2009.^{[5][6]}
- **Internet:** *Google* processes about 20 petabytes of data per day.^[7]
- **Physics:** The 4 experiments in the *Large Hadron Collider* will produce about 15 petabytes of data per year, which will be distributed over the *LHC Computing Grid*.^[8]
- **P2P networks:** As of October 2009, *Isobunt* has about 9.76 petabytes of files contained in *torrents* indexed globally.^[9]
- **Games:** *World of Warcraft* utilizes 1.3 petabytes of storage to maintain its game.^[10]



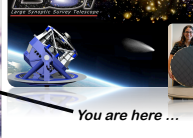
ECS-165B

14

Science has been changing lately ...

- "All science is either physics or stamp collecting."
 - Ernest Rutherford, British chemist & physicist (1871 - 1937)
 - [J. B. Birks "Rutherford at Manchester" (1962)]
- That is, from few data, lots of thinking
- ... to LOTS OF DATA and ANALYSIS

→ "Data-driven" scientific discovery!
4th paradigm, in addition to hypothesis-driven science



You are here ...

ECS-165B

15

Also: Data(bases) can be Yummy!



ECS-165B

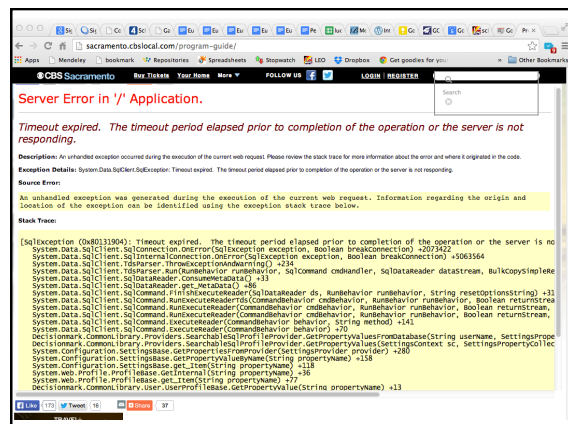
16

Exploits of a Mom <http://xkcd.com/327/>



ECS-165B

17



Introduction to XQuery

Data model

XML

```
<BOOK YEAR="1999 2003">
  <AUTHOR>Abiteboul</AUTHOR>
  <AUTHOR>Buneman</AUTHOR>
  <AUTHOR>Suciu</AUTHOR>
  <TITLE>Data on the Web</TITLE>
  <REVIEW>A <EM>fine</EM> book.</REVIEW>
</BOOK>
```

XQuery

```
element BOOK {
  attribute YEAR { 1999, 2003 },
  element AUTHOR { "Abiteboul" },
  element AUTHOR { "Buneman" },
  element AUTHOR { "Suciu" },
  element TITLE { "Data on the Web" },
  element REVIEW { "A", element EM { "fine" }, "book." }
}
```

XQuery Tutorial

Peter Fankhauser, Fraunhofer IPSI
Peter.Fankhauser@ipsi.fhg.de

Philip Wadler, Avaya Labs
wadler@avaya.com

YOUR USER REQUIREMENTS INCLUDE FOUR HUNDRED FEATURES.

DO YOU REALIZE THAT NO HUMAN COULD BE ABLE TO USE A PRODUCT WITH THAT LEVEL OF COMPLEXITY?

GOOD POINT! TO BEYER ADD "EASY TO USE" TO THE LIST!

XQuery Examples

Titles of all books published before 2000

```
/BOOKS/BOOK[@YEAR < 2000]/TITLE
```

Year and title of all books published before 2000

```
for $book in /BOOKS/BOOK
where $book/@YEAR < 2000
return <BOOK>{ $book/@YEAR, $book/TITLE }</BOOK>
```

Books grouped by author

```
for $author in distinct(/BOOKS/BOOK/AUTHOR) return
  <AUTHOR NAME="{ $author }">{
    /BOOKS/BOOK[AUTHOR = $author]/TITLE
  }</AUTHOR>
```

Handwritten notes:

Input: <author>
Joe Doe
/author>

XPath

Output: <author>
name: Joe Doe
title: FLOWER
title: This My Best
title: How to
title: How to

[<author> ..
;]

MapReduce vs Parallel Databases

MapReduce: A Flexible Data Processing Tool

MapReduce has been used extensively to build a number of applications.

In this article, we introduce the MapReduce programming model, describe the architecture of existing MapReduce systems, and discuss the challenges of scaling MapReduce to large-scale data processing.

MapReduce is a programming model for processing large data sets with a simple abstraction: the user specifies a map function that processes a key-value pair to generate a set of intermediate key-value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. We built a system around this programming model in 2003 to simplify construction of the inverted index for handling searches at Google.com. Since then, more than 10,000 distinct programs have been implemented using MapReduce at Google, including algorithms for large-scale graph processing, text processing, machine learning, and statistical machine translation. The Hadoop open source implementation

MapReduce complements DBMSs since databases are not designed for extract-transform-load tasks, a MapReduce specialty.

MapReduce complements DBMSs since databases are not designed for extract-transform-load tasks, a MapReduce specialty.

MapReduce and Parallel DBMSs: Friends or Foes?

The "no database" (no database) has been hailed as a revolutionary new platform for large-scale, massively parallel data access. Some proponents claim the extreme scalability of MR will replace relational database management systems (RDBMS) in the enterprise. At least one enterprise, Facebook, has implemented a large data warehouse system using MR technology rather than a DBMS.

Here, we argue that using MR systems to perform tasks that are best suited for DBMSs yields less than satisfactory results, concluding that MR is more like an extract-transform-load (ETL) system than a