SE305 Database System Technology

Kenny Q. Zhu Department of Computer Science Shanghai Jiao Tong University

Kenny Zhu



Research Interests:

AI & Knowledge Engineering

Information extraction Knowledge discovery Text Mining & Understanding Natural Language Processing

Domain Specific Languages

Concurrent and Distributed Languages Data Processing Probabilistic Programming

Degrees: National University of Singapore (NUS) Postdoc: Princeton University Experiences: Microsoft Redmond, USA Microsoft Research Asia

Faculty at SJTU since 2009

Advanced Data & Programming Technology Lab (ADAPT-LAB)



PhD:	6
MSc:	8
Undergrads:	Many

3

International Collaboration

















Administrative Info (I)

- All-English Course: everything in English!
- Lecturer:
 - Kenny Zhu, SEIEE #03-541, <u>kzhu@cs.sjtu.edu.cn</u>
 - Office hours: by appointment or after class
- Teaching Assistant: Shanshan Huang
 - Office: SEIEE #-03-341
 - Email: <u>florahuangss@163.com</u>
 - Office hours: Friday 4-5 PM
- Course Web Page (definitive source!):

http://www.cs.sjtu.edu.cn/~kzhu/se305/

Administrative Info (II)

- Textbook: Database System Concepts (5th ed.) By Abraham Silberschatz, Henry F. Korth and S. Sudarshan.
- Lecture materials on course web page (released after the lecture)

Assignments:

- Released on Thursday (usually)
- Due on the following Tuesday (usually)
- Submit hard copies to me/TA during or after class
- Late submission: -30% of full score for each additional day

Administrative Info (III)

3-credit course

Modes of Assessment (tentative):

Quizzes, Assignments:	40%
1 Lab Project (2-3 pax):	30%
Final Exam:	30%

Email me/TA the names and main contact person of your group

Acknowledgement

The slides series used in this course are modified from the slides copyrighted to Abraham Silberschatz, et al. which is available from <u>http://www.db-book.com/</u>.

Course Overview

Introduction

- Relational Model
- SQL
- Database Design and Application Design
- Data Storage
- Query Processing
- Transactions
- Concurrency Control
- Database Architectures
- Parallel and Distributed Databases
- Object-based Databases and XML
- Data Mining and Information Retrieval
- 1-2 Tutorials
- Project Presentations

Introduction

Roadmap For This Lecture

- What is a DBMS?
- Relational Data Models
- Languages to access the database
- Database Design
 - E-R Models
- Object-based data model
- XML
- Storage management
- Transactions
- Database architecture
- History of databases

Database Management System (DBMS)

- DBMS contains information about a particular enterprise
 - Collection of interrelated *data* (the *database*)
 - Set of *programs* to access the data
 - An *environment* that is both *convenient* and *efficient* to use
- Database Applications:
 - Banking: transactions
 - Airlines: reservations, schedules
 - Universities: registration, grades
 - Sales: customers, products, purchases
 - Online retailers: order tracking, customized recommendations
 - Manufacturing: production, inventory, orders, supply chain
 - Human resources: employee records, salaries, tax deductions
- Databases can be very large.
- Databases touch all aspects of our lives

University Database Example

- Application program examples
 - Add new students, instructors, and courses
 - Register students for courses, and generate class rosters
 - Assign grades to students, compute grade point averages (GPA) and generate transcripts
- In the early days, database applications were built directly on top of 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

Drawbacks of using file systems to store data (Cont.)

- 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

Database systems offer solutions to all the above problems

Recent Reverse Trends

- Active research on tackling these challenges in file systems NoSQL movement!
- Map-Reduce, Hadoop, Google Big Tables ...
- Good efficiency for very large data (tera bytes to peta bytes) Database can't store this much data!
- Leverage massive distributed computing resources data centers, server farms, etc.

Levels of Abstraction

Physical level: describes how a record (e.g., instructor) is stored.

Logical level: describes data stored in database, and the relationships among the data.

type *instructor* = **record**

ID : string; *name* : string; *dept_name* : string; *salary* : integer;

end;

View level: application programs hide details of data types. Views can also hide information (such as an employee's salary) for security purposes.

View of Data

An architecture for a database system



Schemas and Instances

- Similar to **types** and **values** in programming languages
- Schema the logical structure of the database
 - Example: The database consists of information about a set of customers and accounts and the relationship between them
 - Analogous to type information of a variable in a program
 - **Physical schema**: database design at the physical level
 - Logical schema: database design at the logical level
 - Changed infrequently
- Instance the actual content of the database at a particular point in time
 - Analogous to the value of a variable
- Physical Data Independence the ability to modify the physical schema without changing the logical schema
 - Applications depend on the logical schema
 - In general, the interfaces between the various levels and components should be well defined so that changes in some parts do not seriously affect others.

Data Models

A collection of tools for describing :

Data , Data relationships, Data semantics, Data constraints

Relational model

• Collection of tables

Entity-Relationship data model (mainly for database design)

Entities (objects) and their inter-relationship

Object-based data models (Object-oriented and Object-relational)

Extension from E-R with encapsulations, methods and object identity

Semi-structured data model (XML)

- Data items of the same type have different attributes
- Other older models:
 - Network model
 - Hierarchical model (Windows registry, XML)

Relational Model



(a) The *instructor* table

A Sample Relational Database

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

dept_name	building	budget
Comp. Sci.	Taylor	100000
Biology	Watson	90000
Elec. Eng.	Taylor	85000
Music	Packard	80000
Finance	Painter	120000
History	Painter	50000
Physics	Watson	70000

(b) The *department* table

Data Manipulation Language (DML)

- Language for accessing and manipulating the data organized by the appropriate data model
 - DML also known as query language
- Two classes of languages
 - Declarative (nonprocedural) user specifies what data is required without specifying how to get those data (e.g. SQL)
 - Procedural user specifies what data is required and how to get those data (e.g. relational algebra, SQL procedural extensions)
- SQL is the most widely used query language
 - Pronounced as "sequel"
 - "Structured Query Language"

Data Definition Language (DDL)

Specification notation for defining the *database schema*

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

- DDL compiler generates a set of table templates stored in a *data dictionary*
 - Data dictionary contains metadata (i.e., data about data)
 - Database schema
 - Integrity constraints
 - Primary key (ID uniquely identifies instructors)
 - Referential integrity (**references** constraint in SQL)
 - e.g. dept_name value in any instructor tuple must appear in department relation
 - Authorization

SQL

SQL: widely used non-procedural language

- Example: Find the name of the instructor with ID 22222
 - select name
 - from instructor
 - where *instructor.ID* = '22222'
- Example: Find the ID and building of instructors in the Physics dept.

select instructor.ID, department.building
from instructor natural join department
where instructor.dept_name = "physics"

Application programs generally access databases through one of

- Language extensions that allow embedded SQL
- Application program interface (e.g., ODBC/JDBC) which allow SQL queries to be sent to a database

Database Design

The process of designing the general structure 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

Database Design?

■ Is there any problem with this design?

ID	name	salary	dept_name	building	budget
22222	Einstein	95000	Physics	Watson	70000
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000
83821	Brandt	92000	Comp. Sci	Taylor	100000
15151	Mozart	40000	Music	Packard	80000
33456	Gold	87000	Physics	Watson	70000
76543	Singh	80000	Finance	Painter	120000

Design Approaches

Normalization Theory

• Formalize what designs are bad, and test for them

Entity Relationship Model

The Entity-Relationship Model

Models an enterprise as a collection of *entities* and *relationships*

- Entity: a "thing" or "object" in the enterprise that is distinguishable from other objects
 - Described by a set of *attributes*
- Relationship: an association among several entities
- Represented diagrammatically by an *entity-relationship diagram:*



What happened to dept_name of instructor and student?

Object-Relational Data Models

Relational model: flat, "atomic" values

Object Relational Data Models

- Extend the relational data model by including object orientation and constructs to deal with added data types.
- Allow attributes of tuples to have complex types, including nonatomic values such as nested relations.
- Preserve relational foundations, in particular the declarative access to data, while extending modeling power.
- Provide backward compatibility with existing relational languages.

XML: Extensible Markup Language

- Semi-structured data model
- Defined by the WWW Consortium (W3C)
- Originally intended as a document markup language not a database language
- The ability to specify new tags, and to create nested tag structures made XML a great way to exchange data, not just documents
- XML has become the basis for all new generation data interchange formats.
- A wide variety of tools is available for parsing, browsing and querying XML documents/data

Storage Management

- Storage manager is a program module that provides the interface between the low-level data stored in the database and the application programs and queries submitted to the system.
- The storage manager is responsible to the following tasks:
 - Interaction with the file manager
 - Efficient storing, retrieving and updating of data
- Issues:
 - Storage access
 - File organization
 - Indexing and hashing

Query Processing

- 1. Parsing and translation
- 2. Optimization
- 3. Evaluation



Query Processing (Cont.)

Alternative ways of evaluating a given query

- Equivalent expressions
- Different algorithms for each operation
- Cost difference between a good and a bad way of evaluating a query can be enormous
- Need to estimate the cost of operations
 - Depends critically on statistical information about relations which the database must maintain
 - Need to estimate statistics for intermediate results to compute cost of complex expressions

Transaction Management

- What if the system fails?
- What if more than one user is concurrently updating the same data?
- A transaction is a collection of operations that performs a single logical function in a database application
- Transaction-management component ensures that the database remains in a consistent (correct) state despite system failures (e.g., power failures and operating system crashes) and transaction failures.
- **Concurrency-control manager** controls the interaction among the concurrent transactions, to ensure the consistency of the database.

Database Users and Administrators



Database System Internals



Database Architecture

The architecture of a database systems is greatly influenced by the underlying computer system on which the database is running:

Centralized

Client-server

- Parallel (multi-processor)
- Distributed

Two-tier and Three-tier Architectures



History of Database Systems

1950s and early 1960s:

- Data processing using magnetic tapes for storage
 - Tapes provided only sequential access
- Punched cards for input
- Late 1960s and 1970s:
 - Hard disks allowed direct access to data
 - Network and hierarchical data models in widespread use
 - Edgar F. Codd defines the relational data model
 - Later won the ACM Turing Award for this work
 - IBM Research begins System R prototype
 - UC Berkeley begins Ingres prototype
 - High-performance (for the era) transaction processing

History (cont.)

1980s:

- Research relational prototypes evolve into commercial systems
 - SQL becomes industrial standard
- Parallel and distributed database systems
- Object-oriented database systems

1990s:

- Large decision support and data-mining applications
- Large multi-terabyte data warehouses
- Emergence of Web commerce
- Early 2000s:
 - XML and XQuery standards
 - Automated database administration
 - Later 2000s:
 - Giant data storage systems (Cloud)
 - ▶ Google BigTable, Yahoo PNuts, Amazon EC2, ...

