

CS209

Computer system design and application

Stéphane Faroult
faroult@sustc.edu.cn

Zhao Yao zhaoy6@sustc.edu.cn

What character streams can do that byte streams cannot

1

They can read (write) lines
when buffered

2

They can change encoding

We have seen last time that Character streams have some capabilities of their own.

3

You can also use a Scanner object with a character stream, which is very good for parsing text input.

You can use a **Scanner**

Very similar to keyboard and screen,
which are character devices

File and Directory Operations

Check the **File** class

Copying, deleting files

Listing and searching directories

and so forth.

You can do a lot of things with files other than reading and writing them. There are constants in the File class that take care of differences between Linux and Windows.

Files USED to be very important.

They still are, to some extents, for specialized applications, and they are still the backbone of persistence. However, as an application developer, you are increasingly isolated from files. I have mentioned the case of Image and Media; I could add Properties, and what we'll see the next times, databases that act as a layer between programs and files. A lot of data comes from networks as well.

Every application, 40 years ago, used to open and close a lot of files, this is no longer the case. You open files mostly to load data in memory, and work there. All this is related to the cost of memory.

1 Mb of memory

1970 \$730,000

When Dennis Ritchie created C, memory was horribly expensive. You had to release memory as soon as you no longer needed it, and try to save bytes.



1 Mb of memory

1990 \$100

Twenty years later, James Gosling could take a more relaxed approach, have a garbage collector freeing memory once in a while, and afford the luxury of a 2-byte char.



Source : www.jcm.it.com



1 Mb of memory

2017

And today? Memory costs next to nothing. Just one problem: as the cost of memory was decreasing, applications were using more and more of it. And computers were supporting more and more users. You'd be wrong to believe that you no longer have to worry about memory. In some languages and environments (I'm thinking of Web servers running PHP) memory-per-user is limited to keep everything under control. But you don't need to fear loading sometimes quite a lot of data in memory.



Source : www.jcm.it.com

Because memory is so cheap these days many people load everything in memory, and save everything when done – or they use databases, which require a different approach.

Load all data at once

Work in memory ➡ Collections

Save everything at once

Use databases when data needs to be shared or for very large volumes

BIG PROBLEM

When working with databases, the logic is **very** different from when you work with files – few Java developers understand it!

The approach that is the correct one with files is the wrong one with databases, because databases are **systems** that are optimized for data retrieval and processing and work much faster than anything you can do in Java.

In **practice**, what are mostly files used for today?

Multimedia, Documents

Handled by specialized programs

Parameters

You have seen it (Media and Image in JavaFx, Property files) in many cases all the file reading is performed by a method in a specialized object and everything is transparent for the developer.

Parameter Files

Specify directories ?

Remote connections?

Java Development Kit

Java Runtime Environment

Parameter files are very important for allowing user to modify some important settings when they just have a .class.

Command-line parameters



javac

Parameter Files

Assign reasonable default values

Replace values found in
parameter file

Replace values specified
on the command line

Locations for what you need
should be in parameter
(property) files.

Program



In **practice**, what are mostly
files used for today?

Multimedia, Documents

Handled by
specialized
programs

Parameters

Data collection and exchange

You mostly read and write data for exchanging data
between systems – which include database systems.

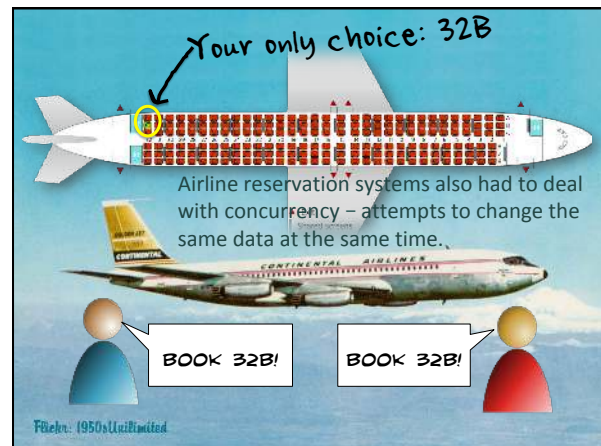
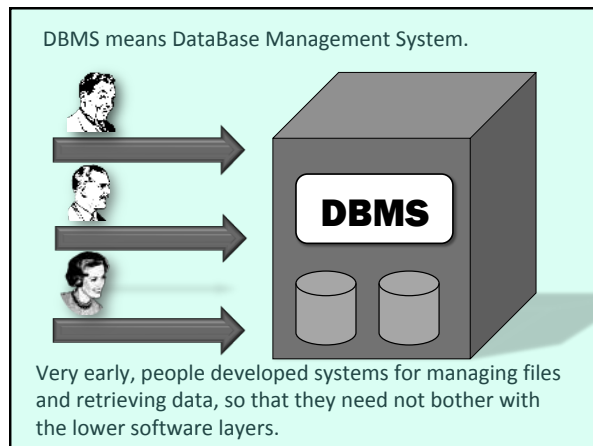
Time to talk about

DATABASES



Computers have always been used for
processing data – which in the 1960s wasn't
as "Big" as today but too big for manual
processing.





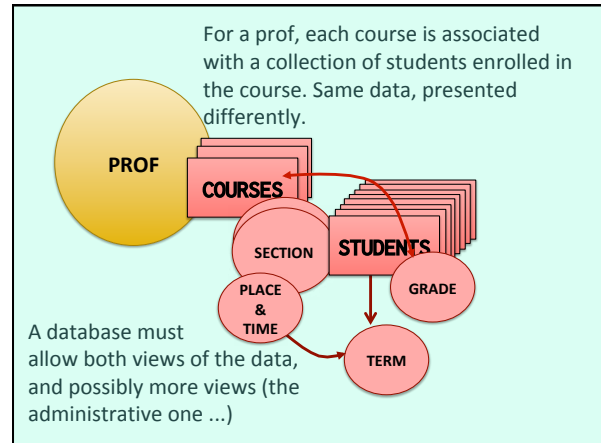
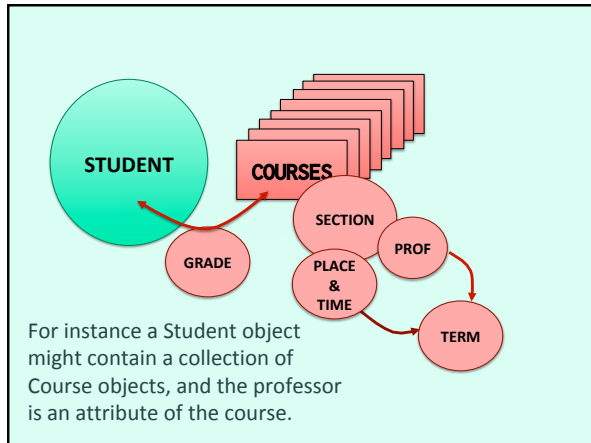
This is what database systems were designed for: checking that changes don't lead to an inconsistent state (two people in the same seat, reservation for a non-existing flight), and also to retrieve data as fast as possible, using a "high-level" language (find this that satisfies those conditions ...) instead of looping on file records and checking each one.

CONTROL changes
RETRIEVE data

Problem of shared data:

DIFFERENT VIEWS

The problem with shared data is that people have different views of the same things. When you design an object for a single application, it's relatively easy.



EARTHQUAKES

id	
UTC date	
latitude	FOX ISLANDS, ALEUTIAN ISLANDS
longitude	ANDREANOF ISLANDS, ALEUTIAN IS.
magnitude	RAT ISLANDS, ALEUTIAN ISLANDS
depth	SOUTH OF ALEUTIAN ISLANDS
region	NEAR ISLANDS, ALEUTIAN ISLANDS

To illustrate the advantages of databases, the region for earthquakes is something hard to process. All these regions refer to the North Pacific area.

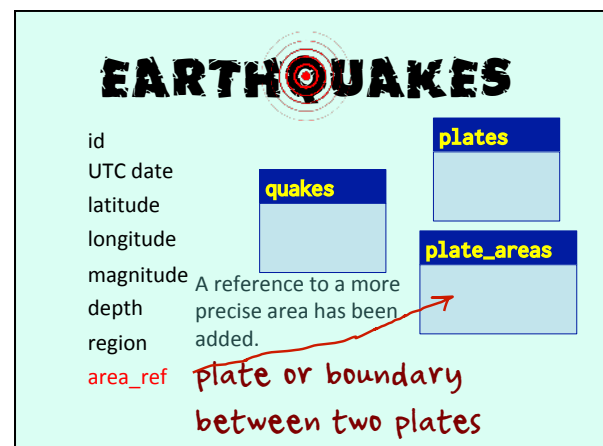
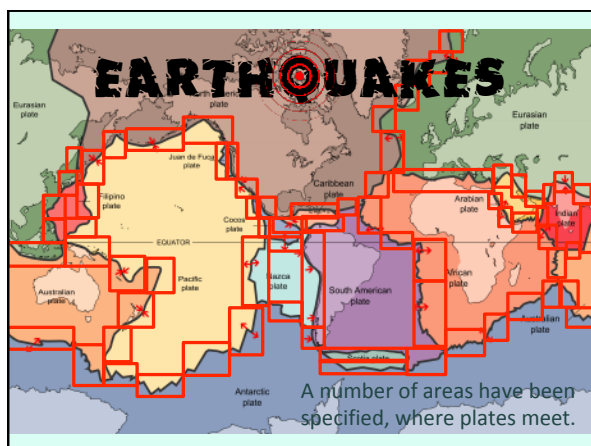
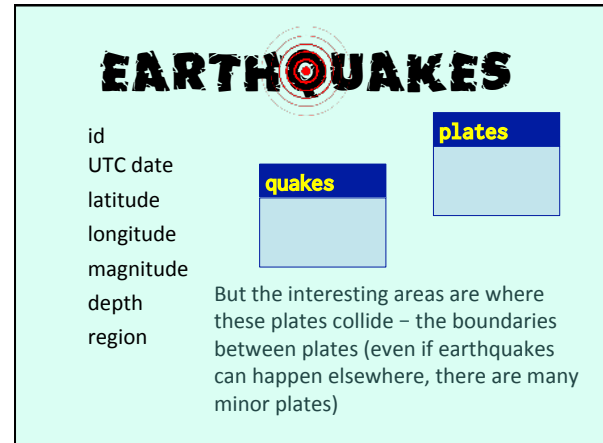
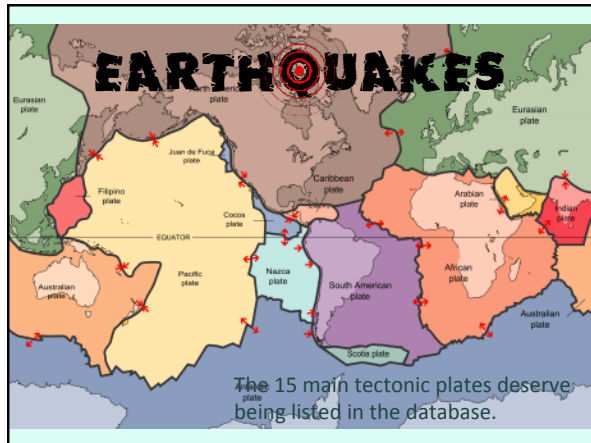
The Fall 2017 project was using an earthquake database (you get one earthquake every 10 minutes on average)

EARTHQUAKES

id	
UTC date	
latitude	
longitude	
magnitude	
depth	
region	

quakes

You can have one collection of quakes but then filtering by region is difficult. So you'd rather add more data to the database, in other tables.



EARTHQUAKES

id
UTC date
latitude
longitude
magnitude
depth
region
area_ref

quakes.csv

plates.csv

plate_areas.csv

?

Now should all this information be stored as .csv files? It would be hard to process.

join would be an entirely separate consideration. In Figure 5, none of the relations R , $\pi_A(R)$, S , $\pi_B(S)$ is a function. Ambiguity in the joining of R with S can sometimes be resolved by means of other relations. Suppose we are given, or can derive from sources independent of R and S , a relation T on the attributes A and B . Then the 3-join of R and S on T is defined by the property that

In 1970, a (British) IBM computer scientist described a way to store data in tables in a database, derived from the mathematical set theory. It was the start of a revolution in data management, and most databases in use today are based on his ideas.

Edgar F. Codd
1923 – 2003

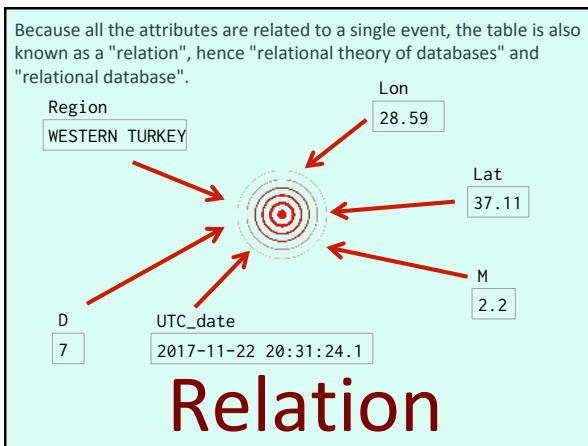
Codd said that the database should give a tabular ("as tables") view of data, which is quite natural. Each table contains a set.

Sets are unordered. If columns are ordered differently, the information isn't changed. Same with rows.

UTC_date	Lat	Lon	D	M	Region
2017-11-22 20:31:24.1	37.11	28.59	7	2.2	WESTERN TURKEY

D	Region	UTC_date	Lon	M	Lat
7	WESTERN TURKEY	2017-11-22 20:31:24.1	28.59	2.2	37.11

But Sets don't allow duplicates. It's important that each different row is stored only once (you define "keys" that allow to specify one particular row. Here, UTC_date, longitude and latitude would be the key)



So ...

It's very tempting to equate what we find in a database to what we use in Object-Oriented programming.

Table = Class

Column = Attribute

Row = Object

?

That would be a very big mistake, unfortunately made by many people.

NO!

In object-oriented programming, the focus is on the object. In a relational database, the focus is on the set.

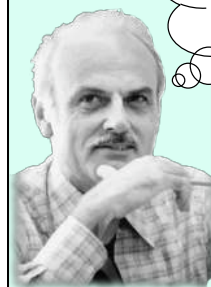
Codd's big idea was that it's possible to operate on sets to define subsets.

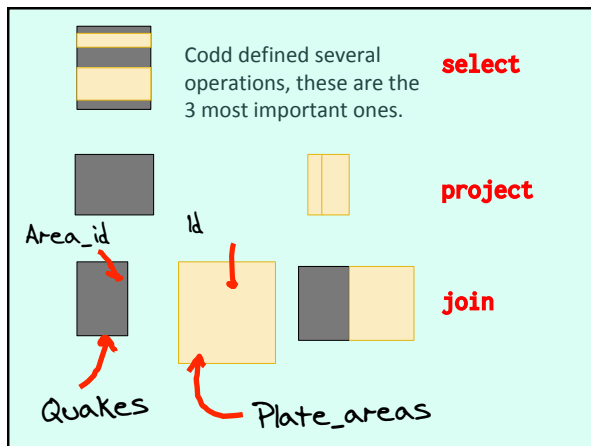
OPERATE
on relations

No references

Link through values

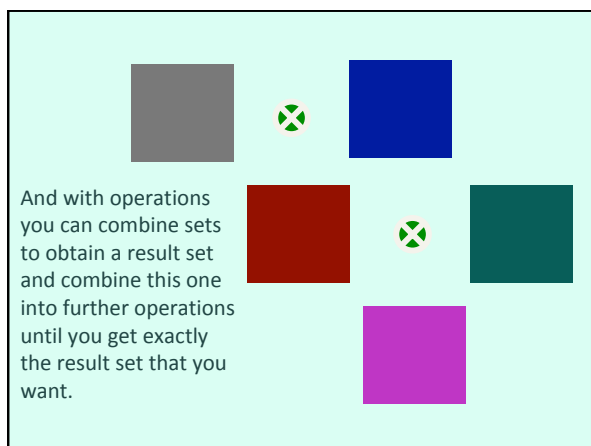
He also insisted on the fact that links between pieces of data in different sets would be established by common data values.





But the great idea, which few people get, is that tables are like variables – instead of containing one value, they contain one set.

Tables
are
variables
(sort of)



One thing that is very important is that tables should be well designed (they must follow a number of rules, called normal forms). You don't want for instance data to be repeated many times, because it would make changes difficult. You must also know precisely what uniquely identifies a row (it may be all columns, but most often it's one or a few columns)

**GOOD MODELLING
IS KEY**

ENTITIES

"Existing things"

student

prof

course

When you design, you look for "entities", things that exist independently of the others (a course can be on the catalogue without any body taking it and teaching it). You must know what identifies each item: a code, a student/employee id, mail address, phone number can be good identifiers. A name isn't good, as several people can have the same name. Then you have attributes. One entity will be one table.

RELATIONSHIPS

Then you have relationships, that link entities together. If an entity can be linked to only one other entity (for instance Student->Dormitory) it can be an attribute of an entity.

student

prof

course

Session

key = combination
of entity keys

Often it will be a table relating identifiers, because a student takes many courses and there are many students in a course.

All this stage of organizing tables is rather difficult and often underrated. If poorly done it can cause many problems.

Distinguishing between what is an attribute and what should be in a relation is often difficult.

NORMALIZATION

What identifies an earthquake?

Date

Latitude

Longitude

Id


I said that date, latitude and longitude identify an earthquake. In the database there is also a sequential number used as identifier. The reason is that you are not supposed to modify an identifier. When an earthquake is first reported, the location isn't always very precise and is sometimes updated later.

SEQUEL: A STRUCTURED ENGLISH QUERY LANGUAGE

by
 Donald D. Chamberlin
 Raymond F. Boyce

Databases are usually associated with a "query language" called SQL, invented in the early 1970s.

IBM Research Laboratory
 San Jose, California



**Don Chamberlin
 with Ray Boyce (+ 1974)**

ABSTRACT: In this paper we present the data manipulation facility for a structured English query language (SEQUEL) which can be used for accessing in an integrated relational data base. Without resorting to the concepts of variables and quantifiers SEQUEL identifies a set of simple operations on tabular structures, which can be shown to be of equivalent power to first order predicate calculus. A SEQUEL user is presented with a constant set of keyword English templates which reflect how people use tables.

SQL was designed to be a very simple language with a syntax close to English that could be used by people who aren't programmers.

**select ...
 from ...
 where ...**

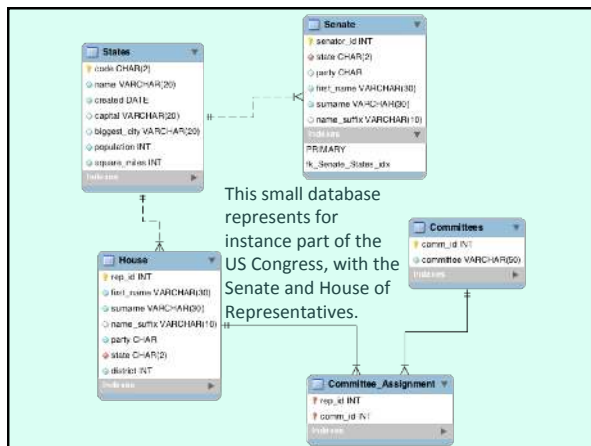
It became at the same time a major success and a major failure: today only computer programmers use SQL – but almost all of them have to use it, and sometimes very often.

TWO main components

SQL provides commands both for managing tables (creating, dropping, modifying them) and for managing data (inserting new rows, updating or deleting existing ones – plus of course commands for retrieving data that satisfies some criteria).

You can learn SQL syntax in about the same time as you can learn how chess pieces move. As with chess, things however quickly become complicated.





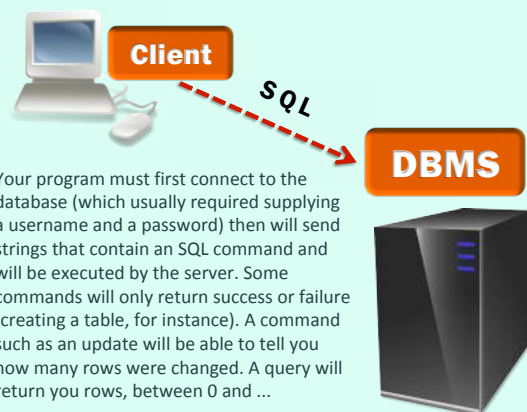
Who participates in the greatest number of House Committees ?

```
select h.first_name, h.surname,
       s.name as state, z.num_of_committees
from (select rep_id,
              count(*) as num_of_committees
      from committee_assignment
      group by rep_id
      having count(*) =
        (select max(num_of_committees)
         from (select rep_id,
                      count(*) as num_of_committees
                  from committee_assignment
                  group by rep_id) x)) z
join house h
  on h.rep_id = z.rep_id
join states s
  on s.code = h.state
order by h.surname, h.first_name, s.name
```

Accessing a database from a Java program is easy, even without using some tools that try to write queries for you (usually they write inefficient queries).

Accessing a database from a Java program

Classic case with a **REAL**
database management
server



Many Java developers have this wrong.

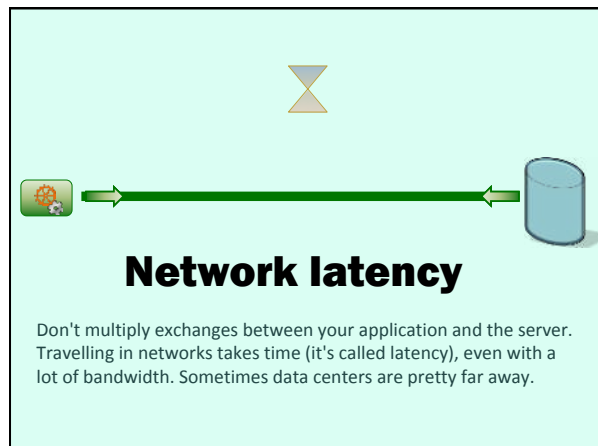
Must SHARE the work between what the database does ...

If you want a sum, don't retrieve all the data and iterate to sum it. Ask for the sum. The DBMS can compute it.

Retrieving
JUST what
you need

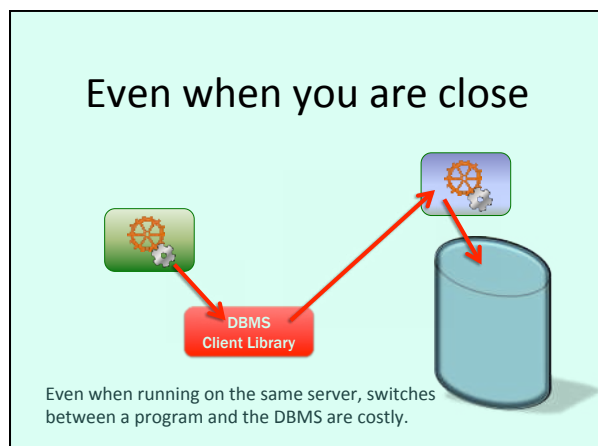
... and what the program does.

Presentation
Computations that
cannot be performed
by the DBMS



If you want to insert many rows into a remote database, send all the rows to the server and ask it to insert all the rows at once (this is called "batching"); it will do it very fast. If you insert rows one by one, you have to wait each time for the response from the server. A return trip between Singapore and Tokyo takes about 0.075s. If you assume that inserting 10,000 rows takes 0.050s (my machine can do that) doing it as one batch between Singapore and Tokyo would take 0.125s. Inserting one row (tested on my machine): 0.003s. Inserting 10000 rows row by row would be $(0.075 + 0.003) * 10000 = 780s$, or 13 minutes ...

ONE insert is OK
10,000 inserts HURT



Embedded databases

No server, single-user

"Connection" same as opening a file

Everything else like the real thing

The alternative to real database servers are "embedded databases", systems that let you use of file as if it were a database.

Apache *Derby*  "Java DB"

Pure Java

Part of the JDK

Server or embedded

Java is shipped with "Derby", sometimes called "Java DB", which can be used both ways.

 SQLite www.sqlite.org

Public domain

One file to download

Used in mobile apps – and by Mozilla

Create tables in *SQLite Manager*
(Firefox Plugin)

The most popular embedded database is probably SQLite.

 SQLite www.sqlite.org



Only for the **C** programming language

<https://github.com/xerial/sqlite-jdbc>

You just need a .jar file available from this address.

sqlite-jdbc refers to JDBC, which is a protocol allowing to access almost all databases from Java. SQL is often slightly different from database to database, not JDBC calls.

jdbc java database connectivity

Load a specific connection driver ^{.jar file}

Database-specific connection parameters

Then java methods are **the same** with every database

SQL is slightly different between databases

Three main Objects

Connection ↔ Link to database, some special operations

PreparedStatement ↔ SQL command

There is also a Statement object. A PreparedStatement can be re-executed with different parameters, not a Statement.

Associate parameters

Execute a query

Loop on rows returned

ResultSet

Database Access Example

Assuming that the database is hosted by an Oracle server, after connection the following program prompts for a date, and issues a query that computes for each region how many quakes happened since that date. The result is displayed line by line.

The toughest part is the connection. Reflection is used for loading the driver, the name of which depends on the DBMS. Parameters required for the connection also depend on the RDBMS.

Java/JDBC

```
import java.util.Properties;
import java.sql.*;
import java.util.Scanner;
```

```
class DBExample {
    static Connection con = null;
```

```
public static void main(String arg[]) {
    Properties info = new Properties();
    String url = "jdbc:oracle:thin:@localhost:1521:orcl";
    Scanner input = new Scanner(System.in);

    try {
        Class.forName("oracle.jdbc.OracleDriver");
    } catch (Exception e) {
        System.err.println("Cannot find the driver.");
        System.exit(1);
    }
}
```

You connect to the database using an url string (not an URL object ...)



Java/JDBC

```
try {
    Class.forName("oracle.jdbc.OracleDriver");
} catch (Exception e) {
    System.err.println("Cannot find the driver.");
    System.exit(1);
}

try {
    System.out.print("Username: ");
    String username = input.nextLine();
    System.out.print("Password: ");
    String password = input.nextLine();
    info.put("user", username);
    info.put("password", password);
    con = DriverManager.getConnection(url, info);
    con.setAutoCommit(false);
    System.out.println("Successfully connected.");
} catch (Exception e) {
    System.err.println(e.getMessage());
    System.exit(1);
}
```

MySQL wants username and password independent parameters, Oracle passes them as Properties. Of course connection can fail.

Java/JDBC

Queries are simple. Placemarkers for parameters in a PreparedStatement are ? signs, implicitly numbered (from 1). Return columns are also numbered from 1. One can often only work with Strings, the DBMS can convert types.

```
System.out.print("Date: ");
String utcDate = input.nextLine();
try {
    PreparedStatement stmt =
        con.prepareStatement("select region, count(*)"
            + " from quakes"
            + " where UTC_date >= ?"
            + " group by region");
    stmt.setString(1, utcDate);
    ResultSet rs = stmt.executeQuery();
    while (rs.next()) {
        System.out.println(rs.getString(1) +
            "\t" + rs.getString(2));
    }
    rs.close();
}
```

Java/JDBC

```
rs.close();
} catch (Exception e) {
    System.err.println(e.getMessage());
    try {
        con.close();
    } catch (SQLException sqLE) {
        // Ignore
    }
    System.exit(1);
}

try {
    con.close();
} catch (SQLException sqLE) {
    // Ignore
}
}
```

Queries can fail too, but returning nothing (or updating or deleting or inserting no rows) doesn't throw any exception, because the empty set is a valid set ... However, trying to insert twice the same key for instance will throw an exception.

```
Class.forName("oracle.jdbc.OracleDriver");  
Class.forName("com.mysql.jdbc.Driver");  
Class.forName("org.postgresql.Driver");  
Class.forName("org.sqlite.JDBC");  
Class.forName("org.apache.derby.jdbc.EmbeddedDriver");
```

These are JDBC driver names for commonly use Database Management systems (note that for Derby there are two modes, embedded – single user – and not embedded). Note also that although most tutorials use reflection to load the driver, you can also use **import** with it, especially with Derby, as the driver will always be in your class path.

And here are connection examples.

```
con = DriverManager.getConnection(url, info);  
    "jdbc:oracle:thin:@hostname:port:dbname"  
    "jdbc:postgresql://hostname:port/dbname"  
    user password  
con = DriverManager.getConnection(url,  
                                username,  
                                password);  
    "jdbc:mysql://hostname:port/dbname"  
con = DriverManager.getConnection(url);  
    "jdbc:sqlite:filename"
```

With Sqlite you just provide a filename. It will be created if it doesn't exist already.