# Introduction to SQL

Group members:

熊国桢、许赞、吴伊宁

曹柳星、王哲、吴宪宇

# What is the hottest website during the National Holiday?

# Dance with Database

# How does it work?

Network Database

Hierarchical Database

→

Relational Database

# Dr. E. F. Codd

《A Relational Model of Data for Large Shared Data Banks》

# Donald D. Chamberlin
# Raymond F. Boyce

**Structured English Query Language**  **Structured Query Language**

| Declarable Programming Language | Procedural Programming Language |
|:---:|:---:|
| **SQL** | **C, C++, Java** |

- Grammar                               Formulation
- Optimization                       More users
- NoSQL                 Huge number of users

# Structured Query Language

- Data Definition

- Data Manipulation

- Data Query

- Data Control

| Category | Verbs | Functions |
| --- | --- | --- |
| **Data Definition Language** | Create, Drop | Create and drop table, view, index |
| **Data Query Language** | Select | Select from tables |
| **Data Manipulation Language** | Insert, Delete, Update | Manipulate data |
| **Data Control Language** | Grant, Revoke | Give or revoke the rights to manipulate a chosen table |

# Grammar of SQL

- **DDL:CREATE&DROP**
- DML:INSERT, DELETE & UPDATE
- DQL: SELECT
- DCL: GRANT&REVOKE

# DDL—CREATE

CREATE TABLE <NAME>
( <Attribute1> Datatype,
<Attribute2> Datatype,…)

Example:
CREATE TABLE user (
USER# INT ,
USERNAME VARCHAR(45) ,
CITY VARCHAR(45) ,
AGE INT ,
PRIMARY KEY (USER#) );

| user |
| --- |
| **USER#** |
| USERNAME |
| CITY |
| AGE |

# DDL—DROP

**DROP TABLE  <NAME>**

Example:
DROP  TABLE test

# Glance at Our Relational DB

**Media**

| |
|---|
| **MEDIA#** |
| MEDIANAME |
| AVGSCORE |
| NO_OF_RATE |
| UP_TIME |
| M_TAG1 |
| M_TAG2 |
| M_TAG3 |

Foreign Key →

**Score**

| |
|---|
| **MEDIA#** |
| SCORE |
| **USER#** |

Foreign Key

Foreign Key

**Collection**

| |
|---|
| **USER#** |
| **MEDIA#** |
| C_TAG1 |
| C_TAG2 |
| C_TAG3 |

Foreign Key

**Users**

| |
|---|
| **USER#** |
| USERNAME |
| CITY |
| AGE |

# Grammar of SQL

- DDL:CREATE&DROP
- **DML:INSERT, DELETE & UPDATE**
- DQL: SELECT
- DCL: GRANT&REVOKE

# DML—INSERT

> **INSERT INTO <NAME>(<Attribute1>,<Attribute2>,...)**
> **VALUES**
> **(<Value of A1>, <Value of A2>,...),**
> **(<Value of A1>, <Value of A2>,...),**
> **...**

Example:

INSERT INTO user(`USER#`,USERNAME,CITY,AGE)
VALUES
(98,   'Penny Liu',   'Shanghai',  19),
(3,    'David Wu',    'Beijing',   21),
(105,  'John Green',  'New York',  18),
(78,   'Amy Liang',   'Xi'an',     23),
(893,  'Suki Hsu',    'Hong Kong', 17);
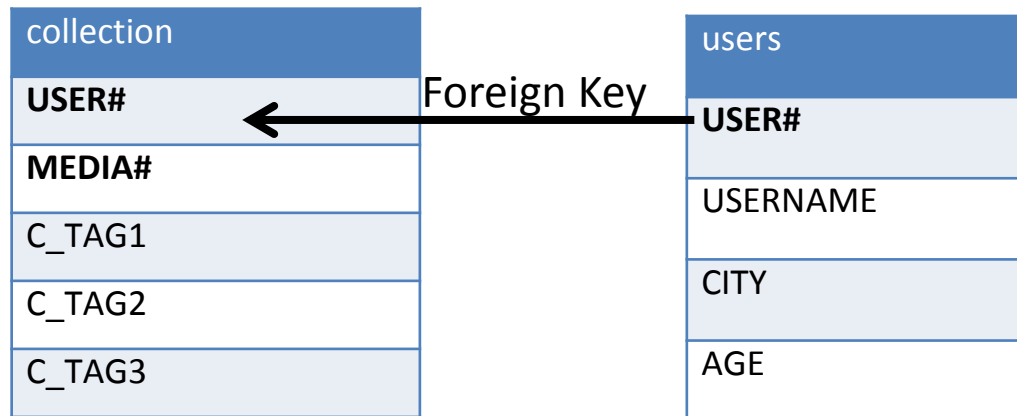
| USER# | USERNAME | CITY | AGE |
|---|---|---|---|
| 98 | Penny Liu | Shanghai | 19 |
| 3 | David Wu | Beijing | 21 |
| 105 | John Green | New York | 18 |
| 78 | Amy Liang | Xi'an | 23 |
| 893 | Suki Hsu | Hong Kong | 17 |

# DML—INSERT(cont'd)

| USER# | USERNAME | CITY | AGE |
|-------|----------|------|-----|
| 98 | Penny Liu | Shanghai | 19 |
| 3 | David Wu | Beijing | 21 |
| 105 | John Green | New York | 18 |
| 78 | Amy Liang | Xi'an | 23 |
| 893 | Suki Hsu | Hong Kong | 17 |

INSERT INTO collection

VALUES

**?**

(135, 1347, `English`, NULL, NULL);

| collection |
|------------|
| **USER#** |
| **MEDIA#** |
| C_TAG1 |
| C_TAG2 |
| C_TAG3 |

Foreign Key ←

| users |
|-------|
| **USER#** |
| USERNAME |
| CITY |
| AGE |

# DML—DELETE

**DELETE FROM<NAME>**
**WHERE<Constraints>**

Example:

DELETE FROM media
WHERE AVGSCORE<6

| MEDIA# | AVGSCORE | NO_OF_RATE |
|--------|----------|------------|
| 1045 | 9.0 | 1 |
| 1734 | 9.0 | 1 |
| 1057 | 7.0 | 1 |
| 1489 | 9.0 | 2 |
| 1347 | 10.0 | 1 |
| 1984 | 5.5 | 4 |
| 1578 | 7.0 | 2 |

| MEDIA# | SCORE | USER# |
|--------|-------|-------|
| 1045 | 9 | 98 |
| 1057 | 7 | 98 |
| 1578 | 8 | 105 |
| 1984 | 3 | 3 |
| 1984 | 6 | 105 |
| 1347 | 10 | 78 |
| 1489 | 8 | 3 |
| 1734 | 9 | 893 |
| 1489 | 10 | 111 |
| 1984 | 7 | 111 |
| 1578 | 6 | 893 |
| 1984 | 6 | 893 |

# DML—UPDATE

> **UPDATE<NAME>**
> **SET<Attribute1>=<Value>**
> **WHERE<Constraints>**

Example:

UPDATE MEDIA

SET

AVGSCORE=(AVGSCORE*NO_OF_RATE+6)/(NO_OF_RATE+1),

NO_OF_RATE=NO_OF_RATE+1

WHERE MEDIA#=1984;

| MEDIA# | AVGSCORE | NO_OF_RATE |
|--------|----------|------------|
| 1045 | 9.0 | 1 |
| 1734 | 9.0 | 1 |
| 1057 | 7.0 | 1 |
| 1489 | 9.0 | 2 |
| 1347 | 10.0 | 1 |
| 1984 | 5.3 | 3 |
| 1578 | 7.0 | 2 |

| MEDIA# | AVGSCORE | NO_OF_RATE |
|--------|----------|------------|
| 1045 | 9.0 | 1 |
| 1734 | 9.0 | 1 |
| 1057 | 7.0 | 1 |
| 1489 | 9.0 | 2 |
| 1347 | 10.0 | 1 |
| 1984 | 5.5 | 4 |
| 1578 | 7.0 | 2 |

# Grammar of SQL

- DDL:CREATE&DROP
- DML:INSERT, DELETE & UPDATE
- **DQL: SELECT**
- DCL: GRANT&REVOKE

# DQL—SELECT

SELECT <Attribute1>,<Attribute2>…
FROM <NAME>
WHERE <Constraints>
Group by <Attribute x>
Having <Constraints>
Order by <attribute>

Example:
SELECT NAME,AGE,Y/N
FROM Tour
WHERE age >= 20 AND Y/N=Y;

## Operation Set

| Tour | Name | Age | Gender | Y/N |
|------|------|-----|--------|-----|
| | Zhang Zheshen | 21 | M | Y |
| | Peng Qijia | 21 | M | Y |
| | Xu Chi | 20 | M | Y |
| | Dai Yibo | 20 | F | Y |
| | Qian Dongliang | 20 | M | Y |
| | Yang Huiqiao | 20 | F | Y |
| | Cao Liuxing | 20 | F | Y |
| | Yuan Quan | 21 | M | Y |

**Selection** $Tour = \sigma_{Y/N=Y \wedge Age \geq 20}(IErs)$

# DQL—SELECT

| USER# | USERNAME | CITY | AGE |
|---|---|---|---|
| 98 | Penny Liu | Shanghai | 19 |
| 3 | David Wu | Beijing | 21 |
| 105 | John Green | New York | 18 |
| 78 | Amy Liang | Xi'an | 23 |
| 893 | Suki Hsu | Hong Kong | 17 |

SELECT USER#,USERNAME,CITY,AGE
FROM user
WHERE CITY = 'Beijing' OR  CITY= 'Xi'an' AND AGE > 22

| USER# | USERNAME | CITY | AGE |
|---|---|---|---|
| 3 | David Wu | Beijing | 21 |
| 78 | Amy Liang | Xi'an | 23 |

What's the result?

SELECT USER#,USERNAME,CITY,AGE
FROM user
WHERE (CITY = 'Beijing' OR  CITY= 'Xi'an') AND AGE > 22

| USER# | USERNAME | CITY | AGE |
|---|---|---|---|
| 78 | Amy Liang | Xi'an | 23 |

24

# DQL—SELECT
## Prep. in WHERE

| Prep. | Example | Meaning |
|---|---|---|
| **NOT**(<Attribute>=<Value>) | **NOT**(CITY=Beijing) | CITY< >Beijing |
| <Attribute> **BETWEEN** <Value1> and <Value2> | AGE **BETWEEN** 18 AND 25 | AGE >=18 **AND** AGE <= 25 |
| <Attribute> **IN**(<Value1>, <Value2>,…) | CITY **IN**(Beijing, Shanghai) | CITY=Beijing **OR** CITY=Shanghai |
| <Attribute> **LIKE** <Description> | CITY **LIKE** 'Bei%' | Value of CITY is like 'Bei%' |

# DQL—SELECT

## Prep. in WHERE(cont'd)

LIKE for Fuzzy Query

-%    represents unknown string with any length>=0

- _    represents exactly **one** character

| WORD BOX | |
|---|---|
| Mp3 | |
| mp3 | |
| song.mp3 | |
| Cambridge | |
| Oxford | |
| mails | |

| A.LIKE  m% | B.LIKE  m_ _ |
|---|---|
| C.LIKE  %m% | D.LIKE  %m_ _ |

# DQL—SELECT
## GROUP BY

| USER# | MEDIA# | C_TAG1 | C_TAG2 | C_TAG3 |
|-------|--------|--------|--------|--------|
| 98 | 1045 | database | model | NULL |
| 893 | 1734 | Japan | Germany | History |
| 893 | 1347 | Car | Spindle | NULL |
| 78 | 1489 | Illumination | NULL | NULL |
| 3 | 1578 | Math | Markov | NULL |
| 78 | 1347 | Spindle | Pro/e | CAM |
| 105 | 1057 | English | Test | TOEFL |

Example:
SELECT
USER#,COUNT(MEDIA#)
FROM collection
GROUP BY USER#
HAVING COUNT(MEDIA#)=1;

| User# | COUNT(Media#) |
|-------|---------------|
| 98 | 1 |
| 3 | 1 |
| 105 | 1 |
| 3 | 1 |
| 105 | 1 |

27

# DQL—SELECT
## Join Tables

| Media |
|---|
| **MEDIA#** |
| MEDIANAME |
| AVGSCORE |
| NO_OF_RATE |
| UP_TIME |
| M_TAG1 |
| M_TAG2 |
| M_TAG3 |

| Collection |
|---|
| **USER#** |
| **MEDIA#** |
| C_TAG1 |
| C_TAG2 |
| C_TAG3 |

| Users |
|---|
| **USER#** |
| USERNAME |
| CITY |
| AGE |



**Operation Set**

**Natural Join**

$$[Detail+] = Detail \bowtie Duration$$

Example:
SELECT c.USER#,u.USERNAME,c.MEDIA#,MEDIANAME,C_tag1,C_tag2,C_tag3
FROM media m,user u,collection c
WHERE m.MEDIA#=c.MEDIA# AND u.USER#=c.USER#

JOIN TABLE ON m.MEDIA#=c.MEDIA#
JOIN TABLE ON u.USER#=c.USER#

28

# DQL—SELECT
Subquery select

| MEDIA# | AVGSCORE | NO_OF_RATE |
|--------|----------|------------|
| 1045 | 9.0 | 1 |
| 1734 | 9.0 | 1 |
| 1057 | 7.0 | 1 |
| 1489 | 9.0 | 2 |
| 1347 | 10.0 | 1 |
| 1984 | 5.5 | 4 |
| 1578 | 7.0 | 2 |

AVG(AVGSCORE)=8.1

SELECT MEDIANAME,MEDIA#,AVGSCORE
FROM media
WHERE AVGSCORE<
(SELECT AVG(AVGSCORE) FROM media)

*Works or Not?*

SELECT MEDIANAME,MEDIA#,AVGSCORE
FROM media
WHERE AVGSCORE< AVG(AVGSCORE)

# DQL—SELECT
## Subquery select

| User# | Media# | C_tag1 | C_tag2 | C_tag3 |
|-------|--------|--------|--------|--------|
| 98 | 1045 | database | model | NULL |
| 893 | 1734 | Japan | Germany | History |
| 893 | 1347 | Car | Spindle | NULL |
| 78 | 1489 | Illumination | NULL | NULL |
| 3 | 1578 | Math | Markov | NULL |
| 78 | 1347 | Spindle | Pro/e | CAM |
| 105 | 1489 | English | Test | TOEFL |

SELECT USER#

FROM collection

WHERE collection.media=**ANY**

(SELECT MEDIA# FROM collection WHERE USER#=78)

AND USER#<>78

## What's the result?

**Answer:893,105**

# Grammar of SQL

- DDL:CREATE&DROP

- DML:INSERT, DELETE & UPDATE

- DQL: SELECT

- **DCL: GRANT&REVOKE**

Example:
GRANT SELECT ON media TO ADMIN1;
REVOKE  ALL      ON media TO ADMIN1;

*(ALL=SELECT,INSERT,DELETE,UPDATE)*

- Grammar                    Formulation
- Optimization               More users
- NoSQL          Huge number of users

# SQL Optimization

- Change the sentence structure
- Using INDEX
- Data mining

# Change the sentence structure

• SNS function: We want to know the number of users who come from Beijing or Tianjin in the database .

- **SELECT** city, **COUNT** (*)

  **FROM** user

  **WHERE** city='Beijing'

  **OR** city='Tianjin'

  **GROUP BY** city

✔ 显示行 0 - 1 ( 2 总计, 查询花费 0.1520 秒)

```
SELECT city, COUNT( * ) AS "NUM"
FROM user
WHERE city = 'Beijing'
OR city = 'Tianjin'
GROUP BY city
```

- **SELECT** city, **COUNT**(*)

  **FROM** user

  **GROUP BY** city

  **HAVING** city = 'Beijing'

  **OR** city = 'Tianjin'

✔ 显示行 0 - 1 ( ~2 总计 ❓, 查询花费 0.1714 秒)

```
SELECT city, COUNT( * ) AS "NUM"
FROM user
GROUP BY city
HAVING city = 'Beijing'
OR city = 'Tianjin'
```

34

# Testing on INDEX

**INDEX**: Data structure that help us retrieve data efficiently and is important when data gets larger.

Step 1: Create a database by phpMyAdmin for testing

| 表 ▲ | 操作 | 行数 ❓ | 类型 | 整理 | 大小 |
|---|---|---|---|---|---|
| collection | 📋 浏览 📝 结构 🔍 搜索 🔧 插入 🗑 清空 ⊝ 删除 | ~309,692 | InnoDB | latin1_swedish_ci | 21 MB |
| media | 📋 浏览 📝 结构 🔍 搜索 🔧 插入 🗑 清空 ⊝ 删除 | ~183,303 | InnoDB | latin1_swedish_ci | 8.5 MB |
| score | 📋 浏览 📝 结构 🔍 搜索 🔧 插入 🗑 清空 ⊝ 删除 | ~314,442 | InnoDB | latin1_swedish_ci | 20 MB |
| user | 📋 浏览 📝 结构 🔍 搜索 🔧 插入 🗑 清空 ⊝ 删除 | ~100,729 | InnoDB | latin1_swedish_ci | 8 MB |
| 4 张表 | 总计 | ~908,166 | InnoDB | latin1_swedish_ci | 57.6 MB |

# Testing on INDEX

Step 2: Test a selection without index

SNS function: we want to find people who live in Beijing and have over 40 years old.



```
1  SELECT userid, username
2  FROM user
3  WHERE age>40
4  AND city = 'Beijing'
5  ORDER BY username
```

Result:  显示行 0 - 29 ( 7,898 总计，查询花费 0.1229 秒)

# Testing on INDEX

Step 3: Create an index in 'username'



| 索引 ⓘ | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 操作 | | | 键名 | 类型 | 唯一 | 紧凑 | 字段 | 基数 | 整理 | 空 | 注释 |
| 🖊 编辑 ⛔ 删除 | | | PRIMARY | BTREE | 是 | 否 | userid | 99907 | A | | 否 | |

| 已用空间 |  |
|---|---|
| 数据 | 5.5 MB |
| 索引 | 0 字节 |
| 总计 | 5.5 MB |

```
1  CREATE INDEX username_idx
2  ON user (username)
```

| 索引 ⓘ | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 操作 | | | 键名 | 类型 | 唯一 | 紧凑 | 字段 | 基数 | 整理 | 空 | 注释 |
| 🖊 编辑 ⛔ 删除 | | | PRIMARY | BTREE | 是 | 否 | userid | 100767 | A | 否 | |
| 🖊 编辑 ⛔ 删除 | | | username_idx | BTREE | 否 | 否 | username | 200 | A | 否 | |

| 已用空间 |  |
|---|---|
| 数据 | 5.5 MB |
| 索引 | 2.5 MB |
| 总计 | 8 MB |

# Testing on INDEX

Step 4: Do the same work as Step 2

在数据库 **xiaoxiong** 运行 SQL 查询：

```
1  SELECT userid, username
2  FROM user
3  WHERE age>40
4  AND city = 'Beijing'
5  ORDER BY username
```

Result: 显示行 0 - 29 ( 7,898 总计，查询花费 0.0034 秒)

⬆

显示行 0 - 29 ( 7,898 总计，查询花费 0.1229 秒)

**40times!!!**

# Summary on INDEX

1. **Effect of INDEX**

显示行 0 - 29 ( 7,898 总计，查询花费 0.0034 秒)

显示行 0 - 29 ( 7,898 总计，查询花费 0.1229 秒)

**40times!!!**

2. **Consideration of other factors**

   Disk space: about 1-1.5 times cost of primary data

   Time spent in updating index

3. **Where needs optimization**

   Create index in frequently-searched attribute

   Create index in attributes which are not often inserted

# Optimization in SQL

Proper INDEX
Good structure of conditions  ⎤ Optimization in time

SQL optimization



Why we need optimization?

# Optimization in Data Mining

Our SNS for learning system has a big quantity of users over 100,000

How could we make progress?

**Data mining**

The process that attempts to discover patterns in large data sets (Wikipedia)
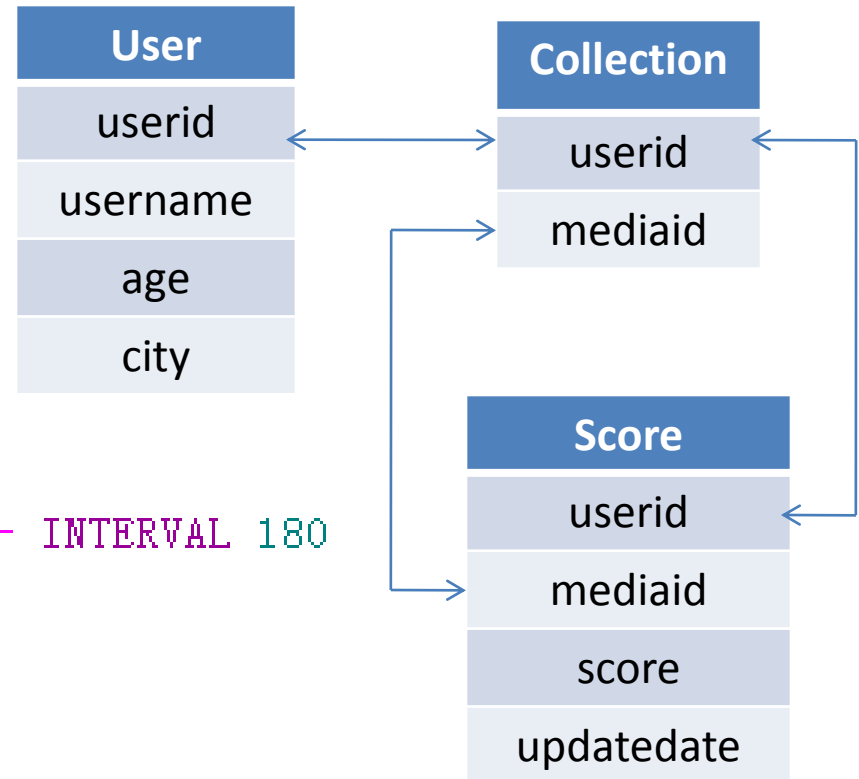
More information of people, more attraction to people

# Optimization in Data Mining

**SNS function**: we wants to search the users who have collected a media object and also scored it in last 180days

There are there tables

| user: | userid | username | city | age |
|---|---|---|---|---|

| collection: | userid | username |
|---|---|---|

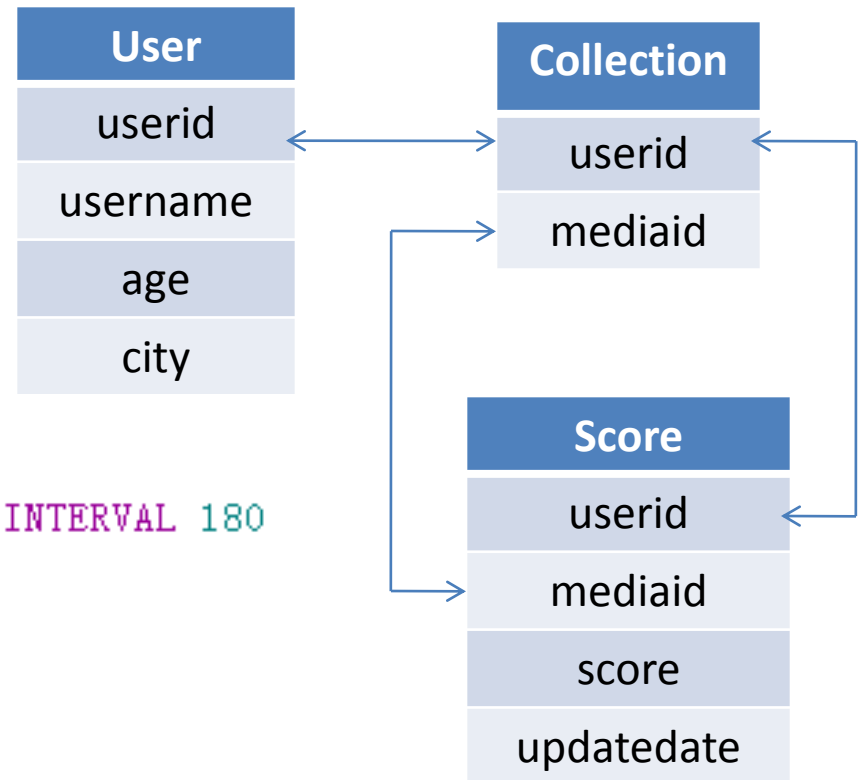| score: | userid | username | score | updatedate |
|---|---|---|---|---|

How to select their '**username**' by using SQL statement.

# Optimization in Data Mining

**SNS function**: we wants to search the users who have collected a media material and also scored it in last 180days

**Is it a solution?**

```
 1 SELECT DISTINCT u.username
 2 FROM user u
 3 JOIN collection c
 4   ON c.userid = u.userid
 5 JOIN score s
 6   ON s.mediaid = c.mediaid
 7 AND s.userid = c.userid
 8 WHERE u.city = 'Beijing'
 9 AND s.updatedate >= CURRENT_Date - INTERVAL 180
10 DAY
```

| User |
|------|
| userid |
| username |
| age |
| city |

| Collection |
|------------|
| userid |
| mediaid |

| Score |
|-------|
| userid |
| mediaid |
| score |
| updatedate |

**Or without "DISTINCT" ?**

# Optimization in Data Mining

Result with DISTINCT

Result without DISTINCT

| username |
|----------|
| pjcnjnrc |
| rxvgdtkc |
| ugfjvcvu |
| cgotnetl |
| riiuocsm |

5 results

| username | userid |
|----------|--------|
| pjcnjnrc | 22672 |
| pjcnjnrc | 22672 |
| rxvgdtkc | 30957 |
| ugfjvcvu | 36633 |
| cgotnetl | 39935 |
| riiuocsm | 60501 |
| riiuocsm | 100001 |

7 results

显示行 0 - 4 ( 5 总计，查询花费 1.0402 秒)

显示行 0 - 6 ( 7 总计，查询花费 1.0396 秒)

Both failed!

44

# Optimization in Data Mining

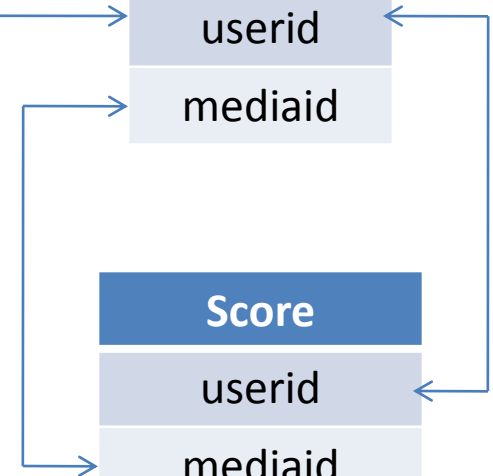**SNS function**: we wants to search the users who have collected a media material and also scored it in last 180days

**Using IN statement**

```
1  SELECT u.username
2  FROM user u
3  WHERE u.userid
4  IN (
5
6  SELECT c.userid
7  FROM collection c, score s
8  WHERE s.mediaid = c.mediaid
9  AND s.userid = c.userid
10 AND s.updatedate >= CURRENT_Date - INTERVAL 180
11 DAY
12 AND u.city = 'Beijing'
13 )
```

| User |
|------|
| userid |
| username |
| age |
| city |

| Collection |
|------------|
| userid |
| mediaid |

| Score |
|-------|
| userid |
| mediaid |
| score |
| updatedate |

# Optimization in Data Mining

Result:



**Certainty   vs.   Time**

# Optimization in Data Mining

**SNS function**: we wants to search the users who have collected a media material and also scored it in last 180days

**Using IN statement**

```
1  SELECT u.username
2  FROM user u
3  WHERE u.userid
4  IN (
5  SELECT c.userid
6  FROM collection c, score s
7  WHERE s.mediaid = c.mediaid
8  AND s.userid = c.userid
9  AND s.updatedate >= CURRENT_Date - INTERVAL 180
10 DAY
11 )
12 AND u.city = 'Beijing'
```

| User |
|------|
| userid |
| username |
| age |
| city |

| Collection |
|------------|
| userid |
| mediaid |

| Score |
|-------|
| userid |
| mediaid |
| score |
| updatedate |

Result: 显示行 0 - 5 ( 6 总计，查询花费 1.1537 秒)

显示行 0 - 5 ( 6 总计，查询花费 4.4028 秒)

# Optimization in Data Mining

**SNS function**: we wants to search the users who have collected a media material and also scored it in last 180days
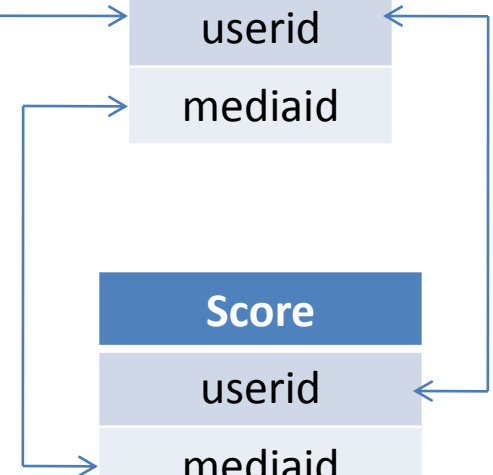
**Using EXISTS statement**

```
1  SELECT DISTINCT u.username
2  FROM user u
3  WHERE EXISTS (
4  SELECT NULL
5  FROM collection c, score s
6  WHERE c.userid = u.userid
7  AND s.mediaid = c.mediaid
8  AND s.userid = c.userid
9  AND s.updatedate >= CURRENT_Date - INTERVAL 180
10 DAY
11 )
12 AND u.city = 'Beijing'
```

Result: 显示行 0 - 5 ( 6 总计，查询花费 1.1209 秒)

Result of IN statement

显示行 0 - 5 ( 6 总计，查询花费 1.1537 秒)

**User**

| userid |
| username |
| age |
| city |

**Collection**

| userid |
| mediaid |

**Score**

| userid |
| mediaid |
| score |
| updatedate |

# Optimization in Data Mining

- **DISTINCT & No DISTINCT**

  Dislodge the repeating data,
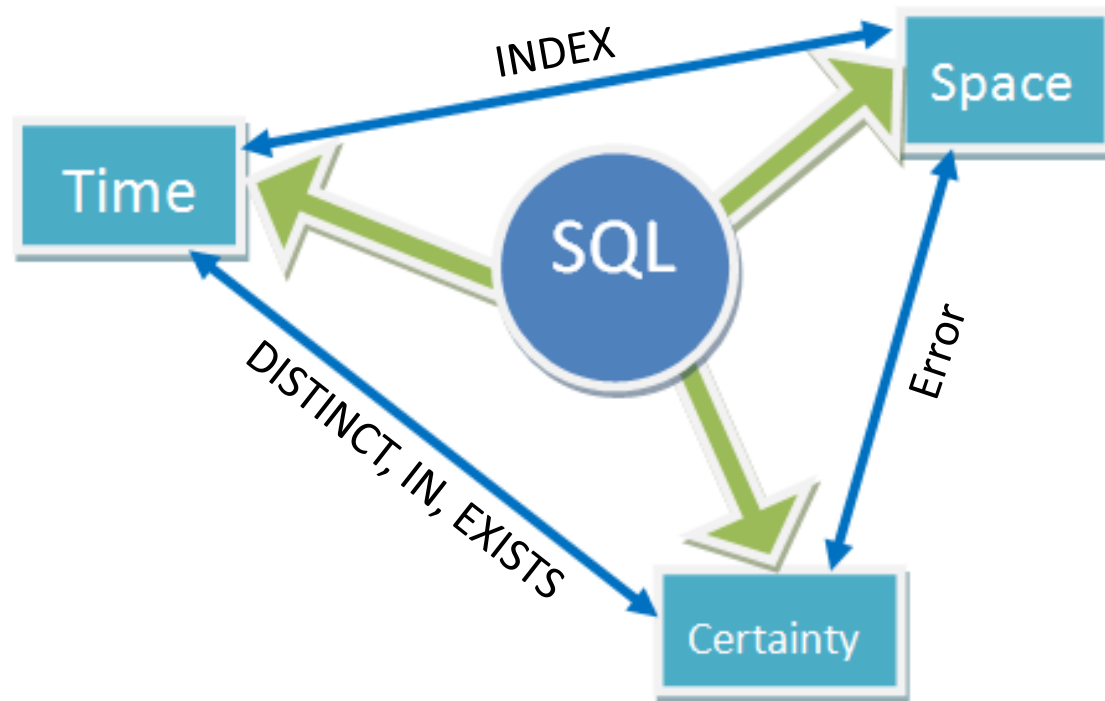  certainty

- **EXISTS & IN**

  Help selection in data mining,
  time, certainty

- **The order of conditions**

  Sharply increase speed of data mining,
  time

# Relation Between SPACE, TIME and CERTAINTY

- Grammar                Formulation
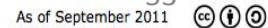- Optimization          More users
- NoSQL          Huge number of users

Can you sure this query language is perfect?

Or the Relational Database is perfect?

Even if the data become huge!!!!!!!!!!!!

As of September 2011

Media
Geographic
Publications
User-generated content
Government
Cross-domain
Life sciences

# Have you heard about these DB?



mongoDB

Cassandra

hadoop

DynamoDB

BigTable

# NoSQL(Not Only SQL)

- **What's meaning of that???**

  NoSQL definition from *nosql-database.org:*

  Next Generation **Databases** mostly addressing some of the points: being **non-relational, distributed, open-source** and **horizontally scalable**.

- **In a word?**

  It does not use SQL as its query language and has a distributed architecture

# Why NoSQL?

- High performance
- High Storage
- High scalability
- High availablity

# What is difference between SQL and NoSQL?

| Database system | Language | Database architerture |
|---|---|---|
| **SQL** | Structured query language(SQL) | Relational |
| **NoSQL** | Java, C, C++, Python, Erlang,etc | **Distributed** |

**Although NoSQL are written in different languages, but they have same architecture**:

**Distributed Database Management System.**

# DDBMS(Distributed Database Management System)

# Popularity of DB    (Oct.2012)

| | | | |
|---|---|---|---:|
| 1. | Oracle | RDBMS | 1619.13 |
| 2. | Microsoft SQL Server | RDBMS | 1242.91 |
| 3. | MySQL | RDBMS | 1232.76 |
| 4. | Microsoft Access | RDBMS | 220.44 |
| 5. | DB2 | RDBMS | 163.52 |
| 6. | PostgreSQL | RDBMS | 144.07 |
| 7. | MongoDB | Document store | 72.92 |
| 8. | SQLite | RDBMS | 64.40 |
| 9. | Cassandra | Wide-column store | 57.96 |
| 10. | Memcached | Key-value store | 21.58 |
| 11. | Redis | Key-value store | 20.90 |
| 12. | HBase | Wide-column store | 15.47 |
| 13. | CouchDB | Document store | 13.97 |
| 14. | Riak | Key-value store | 4.68 |
| 15. | Neo4j | Graph database | 4.33 |
| 16. | Berkeley DB | Key-value store | 2.75 |
| 17. | MariaDB | RDBMS | 2.12 |
| 18. | Oracle NoSQL | Key-value store | 1.72 |

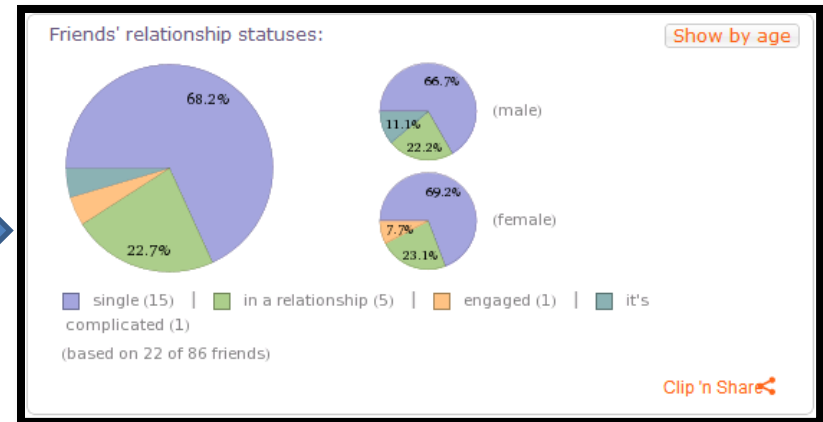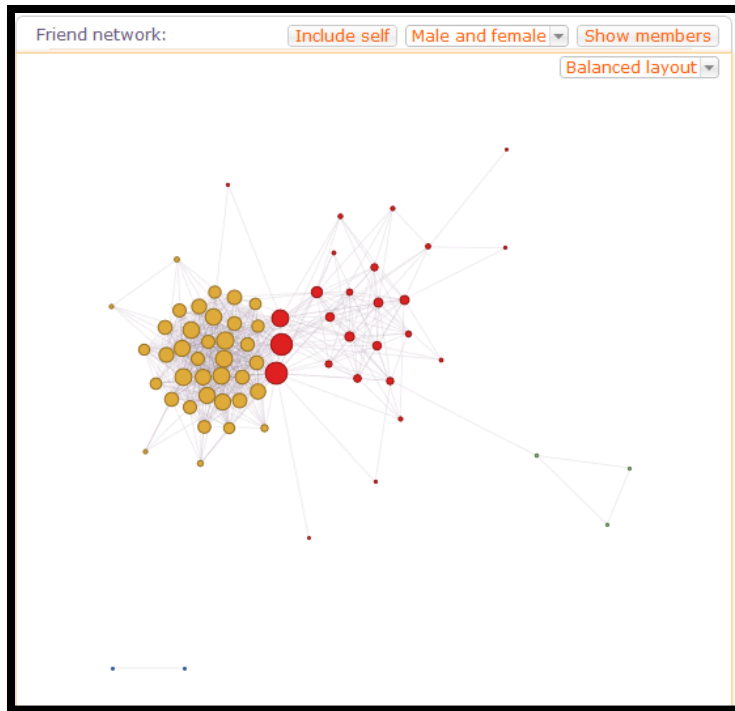# NoSQL is a movement

If we want learn more about these,
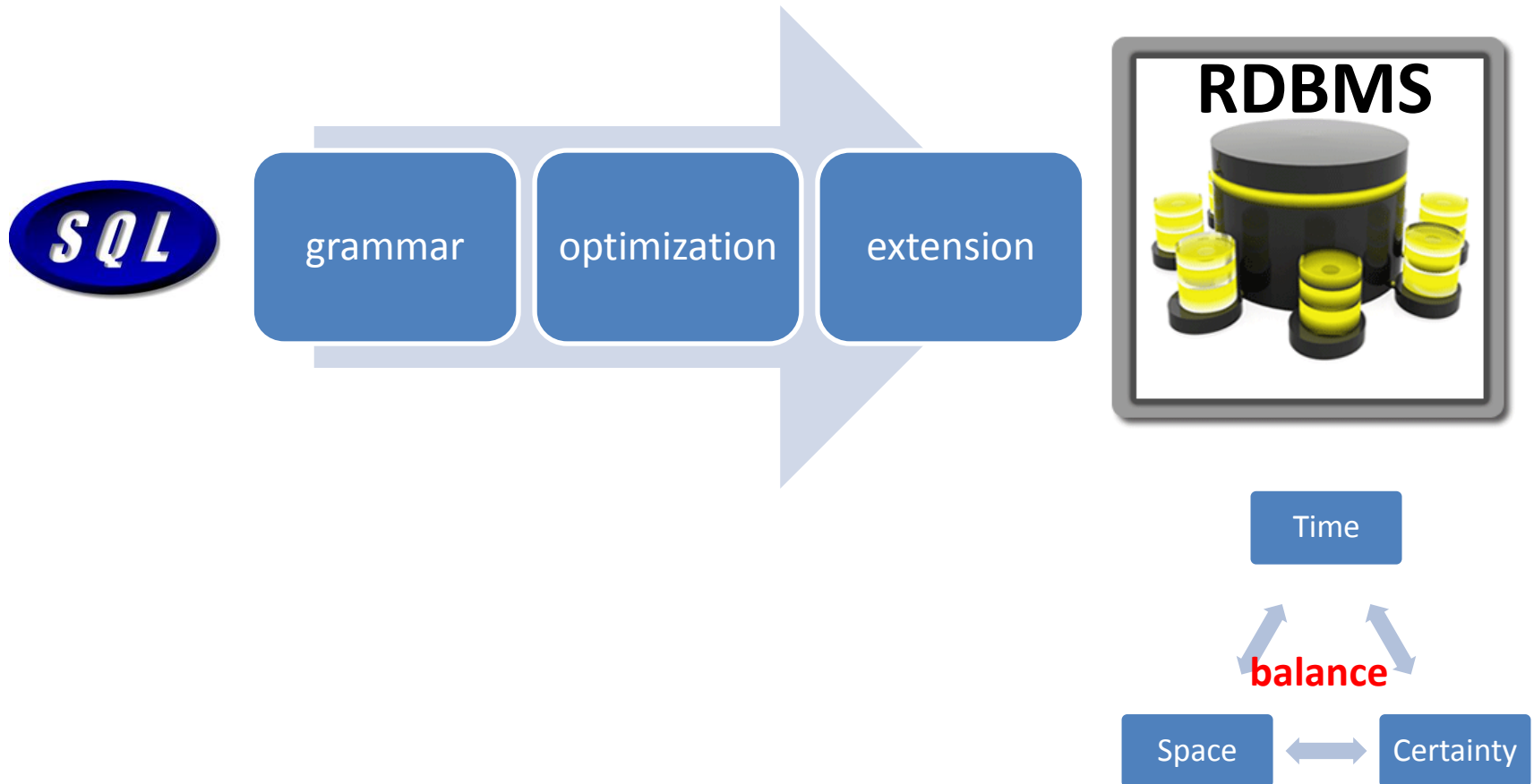
Please learn SQL first!

# Why SQL?

# Why SQL?

- RDBMS, Declarable
- Logically simple, Standardization
- Extension, NoSQL

# SQL for SNS

# Summary

grammar | optimization | extension

**RDBMS**

Time

**balance**

Space ↔ Certainty

# *Refenrence*

- db-engines.com 『*New DB-Engines Ranking shows the popularity of database management systems*』
- http://blogs.x2line.com/al/archive/2007/06/02/3124.aspx 『*Facebook Object-Oriented Diagram*』
- Wikipedia file: 『*LOD Cloud Diagram as of September 2011*』
- 《看图例解学SQL语言》,学苑出版社
- SQL history: *baidu.com*
- RDBMS, DDBMS, NoSQL: *Wikipedia.org*

# Thank you