

Agenda



- Current **trends** in data management & computing
- **Big Data**
- Relational vs. NoSQL databases
 - the value of **relational databases**
 - new **requirements**
 - NoSQL features, strengths and challenges
- **Types** of NoSQL databases
 - **key-value** stores, **document** databases, **column-family** databases, **graph** databases
 - principles and examples

Agenda



- Current trends in data management & computing
- **Big Data**
- Relational vs. NoSQL databases
 - the value of relational databases
 - new requirements
 - NoSQL features, strengths and challenges
- Types of NoSQL databases
 - key-value stores, document databases, column-family databases, graph databases
 - principles and examples

Processing (Traditional) Data



- **OLTP**: Online Transaction Processing
 - Standard **databases** (DBMSs) and database applications
 - Storing, querying, multi-user access
- **OLAP**: Online Analytical Processing (Warehousing)
 - Answer multi-dimensional **analytical** queries
 - Financial/marketing reporting, budgeting, forecasting, ...
- **RTAP**: Real-Time Analytic Processing
(Big Data Architecture & Technology)
 - Data gathered & processed in **real-time** (streaming)
 - Real-time and history **data combined**

Technologies for Big Data



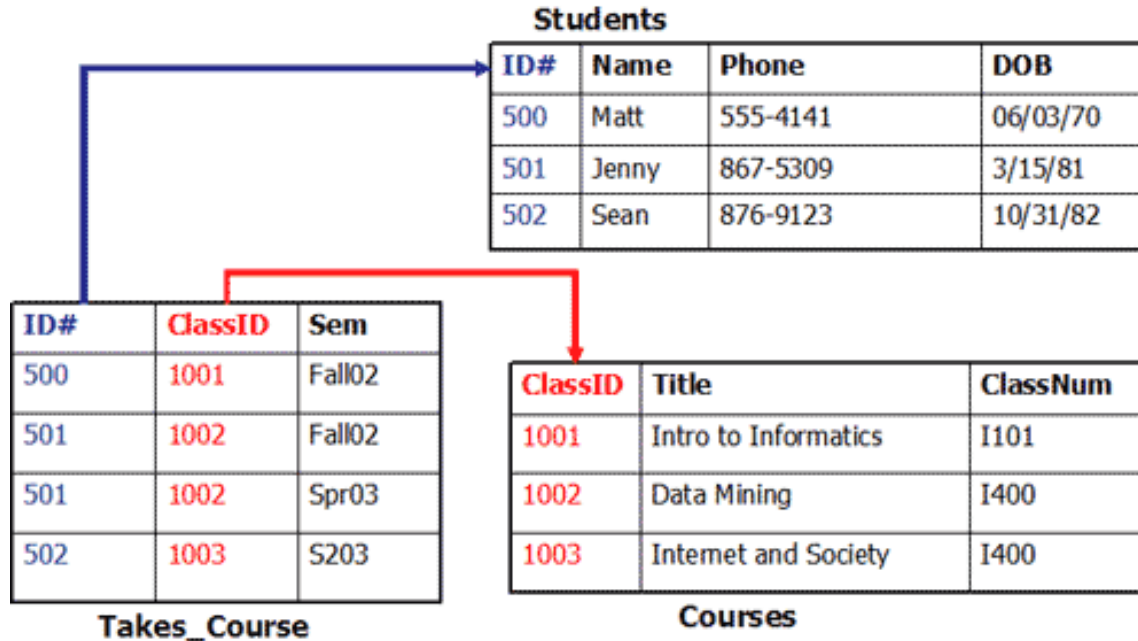
- Distributed file **systems** (GFS, HDFS, etc.)
- **MapReduce**
 - and other models for distributed programming
- **NoSQL databases**
- Data **Warehouses**
- Grid computing, cloud computing
- Large-scale machine learning

Agenda



- Current trends in data management & computing
- Big Data
- Relational vs. NoSQL databases
 - the value of **relational databases**
 - new **requirements**
 - NoSQL features, strengths and challenges
- Types of NoSQL databases
 - key-value stores, document databases, column-family databases, graph databases
 - principles and examples

RDBMS Example



SELECT Name **FROM** Students **NATURAL JOIN**
Takes_Course **WHERE** ClassID = 1001

Data Management: Trends & Requirements



Trends

- **Volume** of data
- **Cloud** comp. (IaaS)
- **Velocity** of data
- **Many** users
- **Variety** of data

Requirements

- Real database **scalability**
 - massive database **distribution**
 - **dynamic** resource management
 - **horizontally** scaling systems
- Frequent **update** operations
- Massive **read** throughput
- **Flexible** database schema
 - semi-structured data

Just Another Temporary Trend?



- There have been **other trends** here before
 - **object** databases, XML databases, etc.
- **But** NoSQL databases:
 - are answer to **real** practical **problems** big companies have
 - are often developed by the **biggest players**
 - outside academia but based on **solid theoretical** results
 - e.g., old results on distributed processing
 - widely used

NoSQL Properties in Detail



1. Good **scalability**

- **horizontal** scalability instead of vertical

2. **Dynamic schema** of data

- different levels of flexibility for **different** types of DB

3. Efficient **reading**

- spend more time to store the data, but **read fast**
- keep relevant information together

4. Cost **saving**

- designed to run on **commodity** hardware
- typically **open-source** (with a support from a company)

Challenges of NoSQL Databases



1. **Maturity** of the technology

- it's getting better, but RDBMS had a lot of time

2. User **support**

- rarely professional support as provided by, e.g. Oracle

3. **Administration**

- massive **distribution** requires advanced administration

4. **Standards** for data access

- RDBMS have SQL, but the NoSQL world is wilder

5. Lack of **experts**

- not enough DB experts on **NoSQL** technologies

Agenda



- Current trends in data management & computing
- Big Data
- Relational vs. NoSQL databases
 - the value of relational databases
 - new requirements
 - NoSQL features, strengths and challenges
- **Types** of NoSQL databases
 - **key-value** stores, **document** databases, **column-family** databases, **graph** databases
 - principles and examples

MapReduce: Features



- MapReduce is a **generic** approach for **distributed** processing of **large** data collections
- **Requires** a way to distribute the **data**
 - and to collect the results back after the processing
- The **user** must only specify two **functions**:
map & reduce

Key-value Stores: Architecture



1. Embedded systems

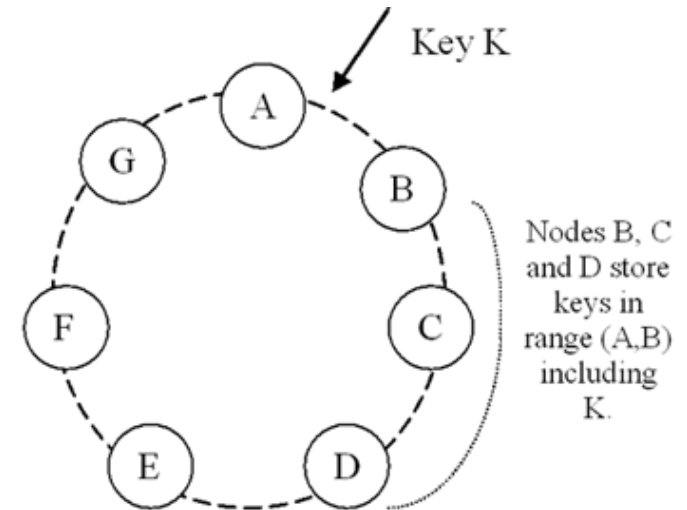
- the system is a **library** and the DB runs **within your** system

2. Large-scale Distributed stores

Architecture often as a **distributed hash table (DHT)**

Features: it is **simple**

- great **performance**, easily scaled



Key-value Stores: Representatives



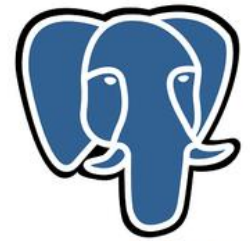
levelDB



redis



Document Databases: Representatives



MS Azure DocumentDB

PostgreSQL

Ranked list: <http://db-engines.com/en/ranking/document+store>

Graph Databases: Mission



- To store **entities** and **relationships** between them
 - **Nodes** are instances of objects
 - Nodes have **properties**, e.g., name
 - **Edges** have **directional** significance
 - Edges have **types** e.g., likes, friend, ...
- Nodes are organized by **relationships**
 - Allow to find interesting patterns
 - example: Get all nodes that are “employee” of “Big Company” and that “likes” “NoSQL Distilled”

Graph Databases: Graphs in RDBMS



- When we store a **graph**-like structure in **RDBMS**, it is for a **single** type of **relationship**
 - “Who is my manager”
- **Adding** another relationship usually means a lot of **schema changes**
- In RDBMS, **we model** the graph **beforehand** based on the **traversal** we want
 - If the traversal changes, the data will have to change
 - **Graph DBs**: the relationship is not calculated but persisted

Facebook: Database Tech. Behind



Apache Hadoop <http://hadoop.apache.org/>



- **Hadoop File System (HDFS)**
 - over 100 PB in a single HDFS cluster
- an open source implementation of **MapReduce**:
 - Enables efficient calculations on massive amounts of data

Apache Hive <http://hive.apache.org/>



- **SQL-like access** to Hadoop-stored data
- integration of **MapReduce** query evaluation

