**CMU SCS**

# Carnegie Mellon Univ.
# Dept. of Computer Science
# 15-415/615 - DB Applications

*C. Faloutsos – A. Pavlo*

Lecture#28: Modern Database Systems
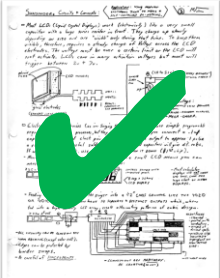
---

**CMU SCS**

# Administrivia – Final Exam

- **Who:** You
- **What:** R&G Chapters 15-22
- **When:** Tuesday May 6th  5:30pm- 8:30pm
- **Where:** WEH 7500
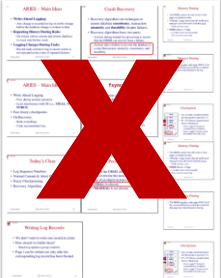- **Why:** Databases will help your love life.

---

**CMU SCS**

# Administrivia – Final Exam



Handwritten Notes              Printed Notes

**CMU SCS**

# Today's Class

- Distributed OLAP
- OldSQL vs. NoSQL vs. NewSQL
- How to scale a database system

**CMU SCS**

# OLTP vs. OLAP

- On-line Transaction Processing:
  - Short-lived txns.
  - Small footprint.
  - Repetitive operations.
- On-line Analytical Processing:
  - Long running queries.
  - Complex joins.
  - Exploratory queries.

**CMU SCS**

# Workload Characterization



Michael Stonebraker – *"Ten Rules For Scalable Performance In Simple Operation' Datastores"*
http://cacm.acm.org/magazines/2011/6/108651

**CMU SCS**

# Relational Database Backlash

- New Internet start-ups hit the limits of single-node DBMSs.
- Early companies used custom middleware to shard databases across multiple DBMSs.
- Google was a pioneer in developing non-relational DBMS architectures.

Faloutsos/Pavlo                    CMU SCS 15-415/615                    7

**CMU SCS**

# MapReduce

- Simplified parallel computing paradigm for large-scale data analysis.
- Originally proposed by Google in 2004.
- Hadoop is the current leading open-source implementation.

Faloutsos/Pavlo                    CMU SCS 15-415/615                    8

**CMU SCS**

# MapReduce Example

Calculate total order amount per day after Jan 1st.

```
REDUCE(key, values) {
    sum = 0;
    while (values.hasNext()) {
        sum += values.next();
    }
    output(key, sum);
}
```

Faloutsos/Pavlo                    CMU SCS 15-415/615                    9

**CMU SCS**

## What MapReduce Does Right

- Since all intermediate results are written to HDFS, if one node crashes the entire query does not need to be restarted.
- Easy to load data and start running queries.
- Great for semi-structured data sets.

Faloutsos/Pavlo          CMU SCS 15-415/615          10

**CMU SCS**

## What MapReduce Did Wrong

- Have to parse/cast values every time:
  - Multi-attribute values handled by user code.
  - If data format changes, code must change.
- Expensive execution:
  - Have to send data to executors.
  - A simple join requires multiple MR jobs.

Faloutsos/Pavlo          CMU SCS 15-415/615          11

**CMU SCS**

## Join Example

- Find sourceIP that generated most adRevenue along with its average pageRank.

12

**CMU SCS**

## Join Example – SQL

```
SELECT INTO Temp sourceIP,
                    AVG(pageRank) AS avgPageRank,
                    SUM(adRevenue) AS totalRevenue
  FROM Rankings AS R, UserVisits AS UV
 WHERE R.pageURL = UV.destURL
   AND UV.visitDate BETWEEN "2000-01-15" AND "2000-01-22"
 GROUP BY UV.sourceIP;

SELECT sourceIP, totalRevenue, avgPageRank
  FROM Temp ORDER BY totalRevenue DESC LIMIT 1;
```

Faloutsos/Pavlo                    CMU SCS 15-415/615                              13

---

**CMU SCS**

## Join Example – MapReduce

| **Phase 1:** Filter | **Phase 2:** Aggregation | **Phase 3:** Search |
|---|---|---|
| **Map:** Emit all records for Rankings. Filter UserVisits data. | **Map:** Emit all tuples (i.e., passthrough) | **Map:** Emit all tuples (i.e., passthrough) |
| **Reduce:** Compute cross product. | **Reduce:** Compute avg pageRank for each sourceIP. | **Reduce:** Scan entire input and emit the record with greatest adRevenue sum. |

Faloutsos/Pavlo                    CMU SCS 15-415/615                              14

---

**CMU SCS**

## Join Example – Results

• Find sourceIP that generated most adRevenue along with its average pageRank.

Chart legend: ☐ Hadoop ■ Vertica ■ DBMS-X

| | 25 nodes | 50 nodes | 100 nodes |
|---|---|---|---|
| Vertica | 32.0 | 35.4 | 55.0 |
| DBMS-X | 29.2 | 29.4 | 31.9 |

15

**CMU SCS**

## Distributed Joins Are Hard

```
SELECT * FROM table1, table2
 WHERE table1.val = table2.val
```

- Assume tables are horizontally partitioned:
  - Table1 Partition Key → table1.key
  - Table2 Partition Key → table2.key
- **Q:** How to execute?
- Naïve solution is to send all partitions to a single node and compute join.

Faloutsos/Pavlo              CMU SCS 15-415/615              16

---

**CMU SCS**

## Semi-Joins

- First distribute the join attributes between nodes and then recreate the full tuples in the final output.
  - Send just enough data from each table to compute which rows to include in output.
- Lots of choices make this problem hard:
  - What to materialize?
  - Which table to send?

Faloutsos/Pavlo              CMU SCS 15-415/615              17

---

**CMU SCS**

## MapReduce in 2014

- SQL/Declarative Query Support
- Table Schemas
- Column-oriented storage.

Faloutsos/Pavlo              CMU SCS 15-415/615              18

**CMU SCS**

# Column Stores

- Store tables as sections of columns of data rather than as rows of data.

---

**CMU SCS**

# Column Stores

```
SELECT sex, AVG(GPA) FROM student
  GROUP BY sex
```

| sid | name | login | age | gpa | sex |
|-----|------|-------|-----|-----|-----|
| 1001 | Faloutsos | christos@cs | 45 | 4.0 | M |
| 1002 | Bieber | jbieber@cs | 21 | 3.9 | M |
| 1003 | Tupac | shakur@cs | 26 | 3.5 | M |
| 1004 | Ke$sha | kesha@cs | 22 | 4.0 | F |
| 1005 | LadyGaGa | gaga@cs | 24 | 3.5 | F |
| 1006 | Obama | obama@cs | 50 | 3.7 | M |

### Row-oriented Storage

<sid,name,login,age,gpa,sex>
<sid,name,login,age,gpa,sex>
<sid,name,login,age,gpa,sex>
<sid,name,login,age,gpa,sex>
<sid,name,login,age,gpa,sex>
<sid,name,login,age,gpa,sex>
<sid,name,login,age,gpa,sex>

---

**CMU SCS**

# Column Stores

```
SELECT sex, AVG(GPA) FROM student
  GROUP BY sex
```

| sid | name | login | age | gpa | sex |
|-----|------|-------|-----|-----|-----|
| 1001 | Faloutsos | christos@cs | 45 | 4.0 | M |
| 1002 | Bieber | jbieber@cs | 21 | 3.9 | M |
| 1003 | Tupac | shakur@cs | 26 | 3.5 | M |
| 1004 | Ke$sha | kesha@cs | 22 | 4.0 | F |
| 1005 | LadyGaGa | gaga@cs | 24 | 3.5 | F |
| 1006 | Obama | obama@cs | 50 | 3.7 | M |

### Column-oriented Storage

sid  name  login  age  gpa  sex

## Column Stores

- Only scan the columns that a query needs.
- Allows for amazing compression ratios:
  - Values for the same query are usually similar.
- Main goal is delay materializing a record back to its row-oriented format for as long as possible inside of the DBMS.
- Inserts/Updates/Deletes are harder…

Faloutsos/Pavlo                CMU SCS 15-415/615                22

## Column Store Systems

- Many column-store DBMSs
  - Examples: Vertica, Sybase IQ, MonetDB
- Hadoop storage library:
  - Example: Parquet, RCFile

Faloutsos/Pavlo                CMU SCS 15-415/615                23

## NoSQL

- In addition to MapReduce, Google created a distributed DBMS called BigTable.
  - It used a GET/PUT API instead of SQL.
  - No support for txns.

- Newer systems have been created that follow BigTable's anti-relational spirit.

Faloutsos/Pavlo                CMU SCS 15-415/615                24

**CMU SCS**

# NoSQL Systems

## Key/Value

redis

riak

FOUNDATIONDB

membase

amazon DynamoDB

## Column-Family

cassandra

APACHE HBASE

HYPERTABLE

## Documents

mongoDB

COUCHBASE

CouchDB relax

RethinkDB

---

**CMU SCS**

# NoSQL Drawbacks

- Developers write code to handle eventually consistent data, lack of transactions, and joins.
- Not all applications can give up strong transactional semantics.

---

**CMU SCS**

# NewSQL

- Next generation of relational DBMSs that can scale like a NoSQL system but without giving up SQL or txns.

## Aslett White Paper

[Systems that] deliver the scalability and flexibility promised by NoSQL while retaining the support for SQL queries and/or ACID, or to improve performance for appropriate workloads.

Matt Aslett – 451 Group (April 4th, 2011)
https://www.451research.com/report-short?entityId=66963

28

## Wikipedia Article

A class of modern relational database systems that provide the same scalable performance of NoSQL systems for OLTP workloads while still maintaining the ACID guarantees of a traditional database system.

Wikipedia (April 2014)
http://en.wikipedia.org/wiki/NewSQL

29

## NewSQL Systems

**New Design**
NUODB
VOLTDB
SAP HANA
deepdb
memsql
Clustrix
Google Spanner

**MySQL Engines**
ScaleDB
Tokutek

**Middleware**
ScaleBase
ScaleArc
dbShards

**CMU SCS**

# NewSQL Implementations

- Distributed Concurrency Control
- Main Memory Storage
- Hybrid Architectures
  - Support OLTP and OLAP in single DBMS.
- Query Code Compilation

Faloutsos/Pavlo CMU SCS 15-415/615 31