# Large-Scale Data Analysis: Bridging the Gap

Alekh Jindal, Yagiz Kargin, Sarath Kumar, Vinay Setty

SAARLAND
UNIVERSITY

COMPUTER SCIENCE
**Information Systems Group**

# Outline

- Motivation: Parallel DBMS vs Map/Reduce
- Schema & Benchmarks Overview
- Original(Pavlo) Map/Reduce Plans
- Improved(SAVY) Design & Implementation
- Improving Hadoop
  - Indexing
  - Co-Partitioning
- Experiments
- Conclusion

# Motivation

- Ever growing data
  - About 20TB per Google crawl!
- Computing Solutions
  - High-end server: 1625.60€/core, 97.66€/GB
  - Share-nothing nodes: 299.50€/core, 166.33 €/GB
- Two Paradigms
  - Parallel DBMS
  - Map/Reduce

# Parallel DBMS



[DeWitt, D. and Gray, J. 1992. ]

# Parallel DBMS: Advantages

- Can be column based
    - Example: Vertica
- Local joins possible
    - Partition based on join key
- Can work on compressed data
    - reduced data transfer
- Flexible query plans
- Supports Declarative languages like SQL

# Parallel DBMS - Shortcomings

- Not free of cost
- Not open source
- Cannot scale to thousands of nodes: why?
  - Less fault tolerant
  - Assumes homogeneous nodes
- Not so easy to achieve high performance
  - Needs highly skilled DBA
  - Needs high maintenance

# Map/Reduce(Hadoop): Advantages

- Free of cost

- Open source

- Fault tolerant

- Scales well to thousands of nodes

- Less maintenance

- Flexible query framework

# Map/Reduce(Hadoop): Shortcomings

- Lack of inbuilt Indexing — Current Focus
- Cannot guarantee local joins — Current Focus
- Performance degradation for SQL like queries — Current Focus
  - Multiple MR phases
  - Each MR phase adds extra cost
- No Flexible query plans
- Data transfer not optimized

# Benchmarks and Schema

# Schema

```sql
CREATE TABLE Documents (
                url VARCHAR
(100) PRIMARY KEY,
                contents TEXT
                );


CREATE TABLE Rankings (
                pageURL VARCHAR
(100) PRIMARY KEY,
                pageRank INT,
                avgDuration INT
```

# Schema

CREATE TABLE **UserVisits** (

                    **sourceIP** VARCHAR(16),

                    **destURL** VARCHAR(100),

                    **visitDate** DATE,

                    **adRevenue** FLOAT,

                    userAgent VARCHAR(64),

                    countryCode VARCHAR(3),

                    languageCode VARCHAR(6),

                    searchWord VARCHAR(32),

                    duration INT

                    );

# Benchmarks 1&2

- **Selection task (Benchmark 1)**
  - SELECT pageURL, pageRank FROM Rankings WHERE pageRank > X;

- **Aggregation task (Benchmark 2)**
  - SELECT sourceIP, SUM(adRevenue) FROM UserVisits GROUP BY sourceIP;
  - SELECT SUBSTR(sourceIP, 1, 7), SUM(adRevenue) FROM UserVisits GROUP BY SUBSTR(sourceIP, 1, 7);

# Benchmark 3: Join Task

Projection & Aggregation

- SELECT INTO Temp **sourceIP, AVG (pageRank) as avgPageRank, SUM (adRevenue)** as totalRevenue

FROM **Rankings** AS R, **UserVisits** AS UV

Join

WHERE **R.pageURL = UV.destURL** AND UV.

visitDate BETWEEN Date('2000-01-15') AND

selection

('2000-01-22') GROUP BY UV.sourceIP;

- SELECT sourceIP, totalRevenue, avgPageRank

FROM Temp **ORDER BY totalRevenue DESC**

**LIMIT 1**;

# Original (Pavlo) MR Plans

# Benchmark 1



SELECT pageURL, pageRank FROM Rankings
WHERE pageRank > 10;

# Benchmark 2: Phase 1

# Benchmark 2: Phase 2

# Benchmark 3 – Phase 1

# Benchmark 3 – Phase 2

# Benchmark 3 – Phase 3

# Improved (Savy) MR Plans

# Binary Data

- Eliminates delimiters
- Avoids splitting
- Makes tuples of fixed length
- Helps in indexing

# Benchmark 3(Design I) – Phase 1

Phase 1

HDFS

Record Reader s
- predicate

Mappers
- Identity

Inter. Resu lt rank

Reducers
- join

<Source IP, URL, PageRank, adReveune>

Result1

User visits | Ranking s

User visits | Ranking s

User visits | Ranking s

Binary data

Easy to classify (just look at record

# Benchmark 3(Design I) – Phase 2

Phase 2

HDFS

Mappers

Combiners

Reducers

<Source IP, URL, PageRank, adReveune>

Result1

<Source IP, URL, PageRank, adReveune>

Result1

<Source IP, URL, PageRank, adReveune>

Result1

Identity

Identity

Identity

Avg(PR), Sum (adRevnue)

Avg(PR), Sum (adRevnue)

Avg(PR), Sum (adRevnu)

Inter. Resu

Inter. Resu

Inter. Resu

Max(Sum (adRevnue)

Source IP, Avg(PR), Sum (adRevenue)

Final Result

No Phase 3!

# Benchmark 3(Design II) – Phase 1

# Benchmark 3(Design I) – Phase

# Improving Hadoop

# Improving Hadoop

- Improve Selection (Indexing)
- Improve Join (Co-partitioning)

# Indexing

- Data Loading
  - index and load data into DFS

- Query Execution
  - index look-up and selection

- Implementation on Hadoop

# Data Loading

- Partitioning
- Sorting
- Bulk Loading
- HID Splits

# Data Loading

# Partitioning

Split input data at tuple boundaries

# Partitioning

Split input data at tuple boundaries

# Partitioning

Split input data at tuple boundaries

# Partitioning

Split input data at tuple boundaries

# Sorting

Sort each split on the index key

| 50 | 23 | 78 | 19 | 3 | 42 |
|----|----|----|----|----|----|

| 60 | 13 | 88 | 17 | 5 | 47 |
|----|----|----|----|----|----|

| 70 | 25 | 57 | 14 | 34 | 45 |
|----|----|----|----|----|----|

# Sorting

Sort each split on the index key

# Bulk Loading

Bulk load CSS tree index

# HID Split

Construct **H**eader-**I**ndex-**D**ata Split

# HID Split

Construct **H**eader-**I**ndex-**D**ata Split

# HID Split

Construct **H**eader-**I**ndex-**D**ata Split



Header:    Index end offset

Data end offset

Start index key

End index key

# HID Split

Construct *H*eader-*I*ndex-*D*ata Split



Header:   Index end offset
          Data end offset
          Start index key
          End index key

# Query Execution

- Partitioning
- Split selection
- Index lookup
- Extractor

# Query Execution

# Partitioning

Read header to get HID boundaries

# Partitioning

Read header to get HID boundaries

# Partitioning

Read header to get HID boundaries

# Partitioning

Read header to get HID boundaries

# Split Selection

Discard splits containing out of range index keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys



Full Contained

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Index Lookup

Find data offsets corresponding to LOW and HIGH keys

# Extractor

Perform selection on data

# Extractor

Pass sub-split to Record Reader for processing

# Implementation on Hadoop

Loading
- CSS Tree Index
- Indirect index
- Four key types supported - Int, Float, Date, String
- Index stored as byte array
- Reducer to reduce number of files
- Integral number of HID splits per reducer output

Querying
- Discover HID split boundaries from respective headers
- Read only the selected data from HDFS

# Co-Partitioning

- Data loading
- Query execution

# Data Loading

# Data Loading

Relation 1

Relation 2

# Data Loading

Relation 1

Relation 2

Group by Join key
(Map)

Group 1 | Group 1

Group 2 | Group 2

Group 3 | Group 3

# Data Loading

# Data Loading

# Query Execution

# Query Execution

# Query Execution

# Query Execution

# Query Execution

# Indexing on top of Co-partitioning

# Indexing on top of Co-partitioning

# Indexing on top of Co-partitioning

# Indexing on top of Co-partitioning

# Indexing on top of Co-partitioning

# Indexing on top of Co-partitioning

# Indexing on top of Co-partitioning

# Experiments

# Experimental Setup

- Hadoop 0.19.1
- 5 nodes
- Speed?
- RAM?
- Gigabit Ethernet
- Data size
  - User Visits: 20GB
  - Rankings: 32MB

# Results



Data Size

# Results

# Results

# Results

# Results

# Roadblocks Faced

- Data generation:
  - 20GB UserVisits, 338MB Rankings in HDFS
  - Took 16 hours for generation
  - Too many OS/library dependencies
  - Poor documentation

- Number of nodes:
  - Allocated 6 nodes
  - Effective (up-and-running) 4 nodes
  - Map/Reduce parallelism not exploited
  - Per-split indexing ideally suited for highly parallel execution

# Roadblocks Faced

- Data normalization
  - Schema uses VARCHAR data types
  - Input data normalized to fixed tuple-sized binaries
  - Byte oriented processing speedup negated by increased input size
  - However, facilitates indexing and co-partitioning

- Low selectivity
  - Selection task has selectivity close to 1
  - Indexing benefits are sabotaged

- Incorrect base result
  - Reported join task result was not correct

# Roadblocks Faced

- Implementation deviation from the paper
  - Composite key is not really used in join task

# Discussion: Loopholes

- Benchmarks are well suited (biased) for databases
- Huge difference in data loading time
- Queries make heavy use of indexing, sorting data
- Query optimization not done for Map/Reduce
- Fault tolerance not compared

# Discussion: We can do better!

- Map/Reduce plans can be optimized
- Normalized binary input data can help
- Indexing feasible and performs good
- Co-partitioning feasible and looks promising

# Conclusions

# References

- Pavlo, A., Paulson, E., Rasin, A., Abadi, D. J., DeWitt, D. J., Madden, S., and Stonebraker, M. 2009. A comparison of approaches to large-scale data analysis. SIGMOD '09.
- DeWitt, D. and Gray, J. 1992. Parallel database systems: the future of high performance database systems. *Commun. ACM* 35, 6 (Jun. 1992)