

A Comparison of Approaches to Large-Scale Data Analysis

Magnus Kirø

Norwegian University of Science and Technology

October 26, 2012

Presentation Goals

The purpose of paper is to consider MapReduce and a regular Database Management Systems for large-scale data analysis.

- 1 **Parallel DBMS** and **MR**, two approaches to large-scale data analysis.
- 2 Thee **Architectural Elements** of DBMS and MR.
- 3 Outline of **Benchmark Results** and **Best Practice** for large-scale data analysis.

1 Two Approaches

- MR vs DBMS
- Schema Support
- Indexing
- Programming Model
- Data Distribution
- Execution Strategy
- Flexibility
- Fault Tolerance

2 Results

- Benchmark Results
- Best Practice

3 Discussion

- Pros'n Cons
- Drawbacks
- Related Topics

4 Conclusion

MR vs DBMS

- MR
 - Only two Functions:
 - **Map**
 - **Reduce**
 - Data is stored in a distributed file system on the node.
- BDMS
 - Tables are partitioned across nodes
 - Query optimizer, that translates SQL to a query plan.
execution of the query plan is divided among multiple nodes.
 - underlying storage details can be disregarded by the programmers.

Schema Support

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Indexing

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Programming Model

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Data Distribution

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Execution Strategy

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Flexibility

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Fault Tolerance

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Benchmark Results

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Best Practice

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

MR vs DBMS

The different setups have different strong points.

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Drawbacks

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Related Topics

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

A more time specific tentative plan

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

This is the last slide.

Any questions?