

Map Reduce

Date: 5/8/2019

Page:

- MR is a Programming model for data Processing
 - MR programs are inherently parallel.
 - Enables massive Scalability across 100/1000 of m/c in Hadoop Cluster.
 - MR - heart of Apache Hadoop.
 - MR refers to 2 separate and distinct task that Hadoop performs.
- 1) Map Job - take set of data and convert into another set of data.
 - Individual elements are broken down into tuples (Key/Value) pair.
 - 2) Reduce Job - take o/p from Map as input and combines those data tuples into smaller set of tuples

Benefits

- 1) Scalability - Can process PB of data stored into HDFS.
- 2) Flexibility - Easier access to mix source of data w/ type of data
- 3) Speed - with parallel processing, minimal data movement Hadoop offer fast processing of massive amount of data
- 4) Simple - Developer can write code in choice of language - Java / C++ / Python.

Date: _____ Page: _____
Block - minimum amount of data that can be read or write

Components with RDBMS

- Seeking is the process of moving object disk's head to a particular place on the disk to read and write data
- characterise the latency of disk op.
- Transfer rate Corresponding to disk's bandwidth.
- Data access dominated by seek - take longer to read or write
- Updating small proportion of record in DB -
BTree works well.
- Updating majority of DB - BTree is less efficient than MR.
MR uses sort / merge to rebuild DB
- MR good fit for Problems
 - Analyze whole data set in batch fashion
 - Particularly adhoc analysis
- TDMS Good fit
 - good point queries or updates where dataset indexed to deliver low latency retrieval
- MR Scales linearly with
 - ↳ Size of data
 - * Data is partitioned
 - * function primitives can work in parallel on separate partitions.
 - * Double the size of input data - job will run extramarks

twice as slowly.

- But if double size of cluster a job will run as fast as the original one.
- not true with SQL queries

<u>features</u>	<u>RDBMS</u>	<u>MR</u>
Data Size	GB	PB
Access	Interactive and Batch	Batch
Updates	Read & Write many times	Write once Read times
Transactions	ACID	None
Structure	Schema on write	Schema on Read
Integrity	High	Low
Scaling	Non linear	Linear

Schema on Read

Interpret data at processing time.

MapReduce

- * Map Reduce is a Programming model for data processing
- * Map Reduce programs are inherently parallel.
- * Distributing task across multiple nodes.
- * Each node process data stored in that nodes
- * Consist of 2 phases - map - reduce

Features Of Map Reduce

- * Automatic parallelization and distribution
- * Fault tolerance
- * Clean abstraction for programmers
- * MR Programs usually written in Java
- * Can be written in any language Using Hadoop Streaming
- * MR developed in Java
- * Abstracts all housekeeping away from developer.
↳ Concentrate on writing MR functions.
- * Data locality
- * In built redundancy.
- * abstract from Complexity of Dist Programming by.
- * manages all inter process communication
- * Shared nothing than architectural model

Cluster - Group of m/c working together to store and process data.

Worker Node

- HDFS to store data - data node
- Map Reduce to process data - task tracker

Master Node

- Name node - manages HDFS
- Job tracker - manages Map Reduce