

Announcements (Wed., Apr. 17)

- Homework #4 due Monday, April 24, 11:55 pm
 No other extra credit problem
- Project Presentations
 A few project groups still have not signed up please sign up soon
- Google Cloud code • Please redeem your code asap, by May 11
- Wednesday April 19
 - Guest Lecture by Prof. Jun Yang
 - OLAP, Data warehousing, Data mining
 - Included in the final









DBMS: The parallel Success Story

- DBMSs are the most successful application of parallelism
 - Teradata (1979), Tandem (1974, later acquired by HP),..
 - Every major DBMS vendor has some parallel server
- Reasons for success:
 - Bulk-processing (= partition parallelism)
 - Natural pipelining
 - Inexpensive hardware can do the trick
 - Users/app-programmers don't need to think in parallel















Example

• R(<u>Key</u>, A, B)

- Can Block-partition be skewed?
 no, uniform
- Can Hash-partition be skewed? • on the key: uniform with a good hash function
 - on A: may be skewed,
 e.g. when all tuples have the same A-value

Parallelizing Sequential Evaluation Code

- "Streams" from different disks or the output of other operators
 - are "merged" as needed as input to some operator
 - are "split" as needed for subsequent parallel processing
- Different Split and merge operations appear in addition to relational operators
- No fixed formula for conversion
- Next: parallelizing individual operations

Parallel Scans

- Scan in parallel, and merge.
- Selection may not require all sites for range or hash partitioning
 - but may lead to skew
 - + Suppose $\sigma_{A\,{\scriptscriptstyle =}\,10}R$ and partitioned according to A
 - Then all tuples in the same partition/processor
- Indexes can be built at each partition

Parallel Sorting



Idea:

- Scan in parallel, and range-partition as you go
 - e.g. salary between 10 to 210, #processors = 20
 - salary in first processor: 10-20, second: 21-30, third: 31-40,
- As tuples come in, begin "local" sorting on each
- Resulting data is sorted, and range-partitioned
- Visit the processors in order to get a full sorted order
- Problem: skew!
- Solution: "sample" the data at start to determine partition points.

Parallel Joins

- Need to send the tuples that will join to the same machine
 - also for GROUP-BY
- Nested loop:
 - Each outer tuple must be compared with each inner tuple that might join
 - Easy for range partitioning on join cols, hard otherwise

Sort-Merge:

- Sorting gives range-partitioning
- · Merging partitioned tables is local



Dataflow Network for parallel Join













Google's MapReduce is inspired by map and reduce functions in functional programming languages. For example, in Scala functional programming language,

scala> val lst = List(1,2,3,4,5)
scala> lst.map(x => x + 1)
reso: List[Int] = List(2,3,4,5,6)

For example, in Scala functional programming language,

scala> val lst = List(1,2,3,4,5)
scala> lst.reduce((a, b) => a + b)
res0: lnt = 15

Ok, it makes sense in one machine.

Then, how does Google extend the functional idea to multiple machines in order to process large data?















Distributed Computation Before MapReduce

Things to consider:

- how to divide the workload among multiple machines?
- how to distribute data and program to other machines?
- how to schedule tasks?
- what happens if a task fails while running?
- ... and ... and ...

Distributed Computation After MapReduce

Things to consider:

- how to write Map function?
- how to write Reduce function?















Fault Tolerance

Although the probability of a machine failure is low, the probability of a machine failing among thousands of machines is common.

How does MapReduce handle machine failures?

Worker Failure

- The master sends heartbeat to each worker node.
- If a worker node fails, the master reschedules the tasks handled by the worker.

Master Failure

• The whole MapReduce job gets restarted through a different master.

Examples

Example problem: Parallel DBMS R(a,b) is horizontally partitioned across N = 3 machines. Each machine locally stores approximately 1/N of the tuples in R. The tuples are randomly organized across machines (i.e., R is block partitioned across machines). Show a RA plan for this query and how it will be executed across the N = 3 machines. Pick an efficient plan that leverages the parallelism as much as possible. • SELECT a, max(b) as topb • FROM R • WHERE a > 0 • GROUP BY a