CSE 444: Database Internals

Lectures 20-21 Parallel DBMSs

CSE 444 - Spring 2015

What We Have Already Learned

- Overall architecture of a DBMS
- Internals of query execution:
 - Data storage and indexing
 - Buffer management
 - Query evaluation including operator algorithms
 - Query optimization
- Internals of transaction processing:
 - Concurrency control: pessimistic and optimistic
 - Transaction recovery: undo, redo, and undo/redo

Where We Are Headed Next

- Scaling the execution of a query (this week)
 - Parallel DBMS
 - Distributed query processing
 - MapReduce
- Scaling transactions (next week)
 - Distributed transactions
 - Replication
- Scaling with NoSQL and NewSQL (in two weeks)

CSE 444 - Spring 2015

Reading Assignments

- Main textbook Chapter 20.1
- Database management systems.
 Ramakrishnan&Gehrke.
 Third Ed. Chapter 22.11

DBMS Deployment: Local



CSE 444 - Spring 2015

DBMS Deployment: Client/Server



DBMS Deployment: 3 Tiers



DBMS Deployment: Cloud







How to Scale?

- We can easily replicate the web servers and the application servers
- We cannot so easily replicate the database servers, because the database is unique
- We need to design ways to scale up the DBMS

How to Scale a DBMS?



What to scale?

OLTP: Transactions per second
 OLTP = Online Transaction Processing

OLAP: Query response time
 OLAP = Online Analytical Processing

Scaling Transactions Per Second

- Amazon
- Facebook
- Twitter
- ... your favorite Internet application...
- Goal is to scale OLTP workloads
- We will get back to this next week

Scaling Single Query Response Time

- Goal is to scale OLAP workloads
- That means the analysis of massive datasets

This Week: Focus on Scaling a Single Query

CSE 444 - Spring 2015

Big Data

• Buzzword?

- Definition from industry:
 - High Volume <u>http://www.gartner.com/newsroom/id/1731916</u>
 - High Variety
 - High Velocity

Big Data

Volume is not an issue

- Databases *do* parallelize easily; techniques available from the 80's
 - Data partitioning
 - Parallel query processing
- SQL is *embarrassingly parallel*
- We will learn how to do this
- And you will implement it in lab 6

Big Data

New workloads are an issue

- Big volumes, small analytics
 - OLAP queries: join + group-by + aggregate
 - Can be handled by today's RDBMSs (e.g., Teradata)
- Big volumes, big analytics
 - More complex Machine Learning, e.g. click prediction, topic modeling, SVM, k-means
 - Requires innovation Active research area

Data Analytics Companies

Explosion of db analytics companies

- Greenplum founded in 2003 acquired by EMC in 2010; A parallel shared-nothing DBMS (this lecture)
- Vertica founded in 2005 and acquired by HP in 2011; A parallel, column-store shared-nothing DBMS
- DATAllegro founded in 2003 acquired by Microsoft in 2008; A parallel, shared-nothing DBMS
- Aster Data Systems founded in 2005 acquired by Teradata in 2011; A parallel, shared-nothing, MapReduce-based data processing system (in two lectures). SQL on top of MapReduce
- Netezza founded in 2000 and acquired by IBM in 2010. A parallel, shared-nothing DBMS.

Great time to be in data management, data mining/statistics, or machine learning!

Two Approaches to Parallel Data Processing

- Parallel databases, developed starting with the 80s (this lecture and next)
 - For both OLTP (transaction processing)
 - And for OLAP (decision support queries)
- MapReduce, first developed by Google, published in 2004 (in two lectures)
 - Only for decision support queries

Today we see convergence of the two approaches (Greenplum, BigQuery)

Parallel DBMSs

- Goal
 - Improve performance by executing multiple operations in parallel
- Key benefit
 - Cheaper to scale than relying on a single increasingly more powerful processor
- Key challenge
 - Ensure overhead and contention do not kill performance

Performance Metrics for Parallel DBMSs

Speedup

- More processors → higher speed
- Individual queries should run faster
- Should do more transactions per second (TPS)
- Fixed problem size overall, vary # of processors ("strong scaling")

Linear v.s. Non-linear Speedup



Performance Metrics for Parallel DBMSs

Scaleup

- More processors → can process more data
- Fixed problem size per processor, vary # of processors ("weak scaling")
- Batch scaleup
 - Same query on larger input data should take the same time
- Transaction scaleup
 - N-times as many TPS on N-times larger database
 - But each transaction typically remains small

Linear v.s. Non-linear Scaleup



Warning

- Be careful. Commonly used terms today:
 - "scale up" = use an increasingly more powerful server
 - "scale out" = use a larger number of servers

Challenges to Linear Speedup and Scaleup

• Startup cost

Cost of starting an operation on many processors

Interference

Contention for resources between processors

- Skew
 - Slowest processor becomes the bottleneck

Three Architectures for Parallel DB

- Shared memory
- Shared disk
- Shared nothing

Architectures for Parallel Databases

Figure 1 - Types of database architecture





From: Greenplum Database Whitepaper

CSE 444 - Spring 2015

Shared Memory

- Nodes share both RAM and disk
- Dozens to hundreds of processors

Example: SQL Server runs on a single machine and can leverage many threads to get a query to run faster (see query plans)

- Easy to use and program
- But very expensive to scale

Shared Disk

- All nodes access the same disks
- Found in the largest "single-box" (non-cluster) multiprocessors

Oracle dominates this class of systems

Characteristics:

 Also hard to scale past a certain point: existing deployments typically have fewer than 10 machines

Shared Nothing

- Cluster of machines on high-speed network
- Called "clusters" or "blade servers"
- Each machine has its own memory and disk: lowest contention.

NOTE: Because all machines today have many cores and many disks, then shared-nothing systems typically run many "nodes" on a single physical machine.

Characteristics:

- Today, this is the most scalable architecture.
- Most difficult to administer and tune.

We discuss only Shared Nothing in class

In Class

- You have a parallel machine. Now what?
- How do you speed up your DBMS?

Taxonomy for Parallel Query Evaluation

- Inter-query parallelism
 - Each query runs on one processor



Taxonomy for Parallel Query Evaluation

- Inter-query parallelism
 - Each query runs on one processor
- Inter-operator parallelism
 - A query runs on multiple processors
 - An operator runs on one processor




Taxonomy for Parallel Query Evaluation

- Inter-query parallelism
 - Each query runs on one processor
- Inter-operator parallelism
 - A query runs on multiple processors
 - An operator runs on one processor
- Intra-operator parallelism
 - An operator runs on multiple processors



Product Purchase

M

cid=cid

M id=cid

Taxonomy for Parallel Query Evaluation

- Inter-query parallelism
 - Each query runs on one processor
- Inter-operator parallelism
 - A query runs on multiple processors
 - An operator runs on one processor
- Intra-operator parallelism
 - An operator runs on multiple processors



M

cid=cid

cid=cid

M id=cid

We study only intra-operator parallelism: most scalable

Parallel Query Processing

How do we compute these operations on a shared-nothing parallel db?

- Selection: $\sigma_{A=123}(R)$
- Group-by: $\gamma_{A,sum(B)}(R)$
- Join: R [⋈] S

Before we answer that: how do we store R (and S) on a shared-nothing parallel db?







- Relation R split into P chunks R₀, ..., R_{P-1}, stored at the P nodes
- Block partitioned
 - Each group of k tuples goes to a different node
- Hash based partitioning on attribute A:
 - Tuple t to chunk h(t.A) mod P
- Range based partitioning on attribute A:

- Tuple t to chunk i if $v_{i-1} < t.A < v_i$

- Let R(K,A,B,C); which of the following partition methods may result in skewed partitions?
- Block partition
- Hash-partition
 - On the key K
 - On the attribute A
- Range-partition
 - On the key K
 - On the attribute A

Uniform

Uniform

Assuming uniform

hash function

- Let R(K,A,B,C); which of the following partition methods may result in skewed partitions?
- Block partition
- Hash-partition
 - On the key K
 - On the attribute A
- Range-partition
 - On the key K
 - On the attribute A

 Let R(K,A,B,C); which of the following partition methods may result in skewed partitions?



 Let R(K,A,B,C); which of the following partition methods may result in skewed partitions?



CSE 444 - Spring 2015

Data Partitioning Revisited

What are the pros and cons?

- Block based partitioning
 - Good load balance but always needs to read all the data
- Hash based partitioning
 - Good load balance
 - Can avoid reading all the data for equality selections
- Range based partitioning
 - Can suffer from skew (i.e., load imbalances)
 - Can help reduce skew by creating uneven partitions

All three choices are just special cases:

- For each tuple, compute bin = f(t)
- Different properties of the function *f* determine hash vs. range vs. round robin vs. anything

Parallel Selection

Compute $\sigma_{A=v}(R)$, or $\sigma_{v1 < A < v2}(R)$

- On a conventional database: cost = B(R)
- Q: What is the cost on a parallel database with P processors ?
 - Block partitioned
 - Hash partitioned
 - Range partitioned

Parallel Selection

Compute $\sigma_{A=v}(R)$, or $\sigma_{v1 < A < v2}(R)$

- On a conventional database: cost = B(R)
- Q: What is the cost on a parallel database with P processors ? A: B(R) / P, but

 - Hash partitioned
 - Range partitioned
 - Block partitioned
 -- all servers do the work
 - -- one server does the work
 - -- some servers do the work

Data: R(K,A,B,C) -- hash-partitioned on K Query: $\gamma_{A,sum(B)}(R)$



CSE 444 - Spring 2015

- Step 1: each server i partitions its chunk R_i using a hash function h(t.A) mod P: R_{i,0}, R_{i,1}, ..., R_{i,P-1}
- Step 2: server j computes $\gamma_{A, sum(B)}$ on $R_{0,j}, R_{1,j}, ..., R_{P-1,j}$

Compute $\gamma_{A,sum(B)}(R)$

- On a conventional database: cost = B(R)
- Q: What is the cost on a parallel database with P processors ?

Compute $\gamma_{A,sum(B)}(R)$

- On a conventional database: cost = B(R)
- Q: What is the cost on a parallel database with P processors ?
- A: B(R) / P

Can we do better?

- Sum?
- Count?
- Avg?
- Max?
- Median?

Can we do better?

- Sum?
- Count?
- Avg?
- Max?
- Median?

Distributive	Algebraic	Holistic
$sum(a_{1}+a_{2}++a_{9})=sum(sum(a_{1}+a_{2}+a_{3})+sum(a_{4}+a_{5}+a_{6})+sum(a_{7}+a_{8}+a_{9}))$	avg(B) = sum(B)/count(B)	median(B)

- Data: R(<u>K1</u>,A, C), S(<u>K2</u>, B, D)
- Query: R(<u>K1</u>,A,C) ⋈ S(<u>K2</u>,B,D)

Initially, both R and S are horizontally partitioned on K1 and K2







- Data: R(K1,A, C), S(K2, B, D)
- Query: $R(K1,A,C) \bowtie S(K2,B,D)$

Initially, both R and S are horizontally partitioned on K1 and K2



- Step 1
 - Every server holding any chunk of R partitions its chunk using a hash function h(t.A) mod P
 - Every server holding any chunk of S partitions its chunk using a hash function h(t.B) mod P
- Step 2:
 - Each server computes the join of its local fragment of R with its local fragment of S

Compute R $\bowtie_{A=B}$ S

- On a conventional database: cost = B(R)+B(S)
- Q: What is the cost on a parallel database with P processors ?

Compute $R \bowtie_{A=B} S$

- On a conventional database: cost = B(R)+B(S)
- Q: What is the cost on a parallel database with P processors ?
- A: (B(R)+B(S)) / P

Speedup and Scaleup

- Consider:
 - Query: $\gamma_{A,sum(C)}(R)$
 - Runtime: dominated by reading chunks from disk
- If we double the number of nodes P, what is the new running time?
- If we double both P and the size of R, what is the new running time?

Speedup and Scaleup

- Consider:
 - Query: $\gamma_{A,sum(C)}(R)$
 - Runtime: dominated by reading chunks from disk
- If we double the number of nodes P, what is the new running time?
 - Half (each server holds ½ as many chunks)
- If we double both P and the size of R, what is the new running time?
 - Same (each server holds the same # of chunks)

Optimization for Small Relations

When joining R and S

- If |R| >> |S|
 - Leave R where it is
 - Replicate entire S relation across nodes
- Also called a small join or a broadcast join

Other Interesting Parallel Join Implementation

Skew:

- Some partitions get more input tuples than others Reasons:
 - Range-partition instead of hash
 - Some values are very popular:
 - Heavy hitters values; e.g. 'Justin Bieber'
 - Selection before join with different selectivities
- Some partitions generate more output tuples than others

Some Skew Handling Techniques

If using range partition:

- Ensure each range gets same number of tuples
- E.g.: $\{1, 1, 1, 2, 3, 4, 5, 6\} \rightarrow [1,2]$ and [3,6]
- Eq-depth v.s. eq-width histograms

Some Skew Handling Techniques

Create more partitions than nodes

- And be smart about scheduling the partitions
- Note: MapReduce uses this technique

Some Skew Handling Techniques

Use subset-replicate (a.k.a. "skewedJoin")

- Given R ⋈_{A=B} S
- Given a heavy hitter value R.A = 'v' (i.e. 'v' occurs very many times in R)
- Partition R tuples with value 'v' across all nodes e.g. block-partition, or hash on other attributes
- Replicate S tuples with value 'v' to all nodes
- R = the build relation
- S = the probe relation

Parallel Query Evaluation

- Parallel query plan: tree of parallel operators Intra-operator parallelism
 - Data streams from one operator to the next
 - Typically all cluster nodes process all operators
- Can run multiple queries at the same time Inter-query parallelism
 - Queries will share the nodes in the cluster

Parallel Query Evaluation

New operator: Shuffle

- Handles data routing, buffering, and flow control
- Inserted between consecutive operators in the query plan
- Two components: ShuffleProducer and ShuffleConsumer
- Producer:
 - Pulls data from operator and sends to n consumers
 - Producer acts as driver for operators below it in query plan
- Consumer:
 - Buffers input data from n producers and makes it available to operator through getNext interface

Example: Teradata – Loading



AMP = "Access Module Processor" = unit of parallelism

CSE 444 - Spring 2015
Example: Teradata – Query Execution

Find all orders from today, along with the items ordered



Order(oid, item, date), Line(item, ...)







CSE 444 - Spring 2015

Order(oid, item, date), Line(item, ...)







CSE 444 - Spring 2015

Order(oid, item, date), Line(item, ...)

Query Execution



CSE 444 - Spring 2015