

## CSE 444: Database Internals

Lectures 20-21  
Parallel DBMSs

CSE 444 - Spring 2016

1

## 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

CSE 444 - Spring 2016

2

## Where We Are Headed Next

- Scaling the execution of a query
  - Parallel DBMS
  - Distributed query processing
  - MapReduce
- Scaling transactions
  - Distributed transactions
  - Replication
- Scaling with NoSQL and NewSQL

CSE 444 - Spring 2016

3

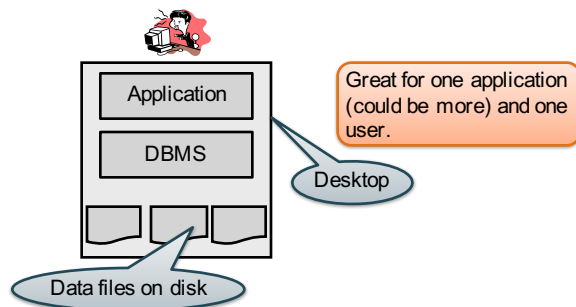
## Reading Assignments

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

CSE 444 - Spring 2016

4

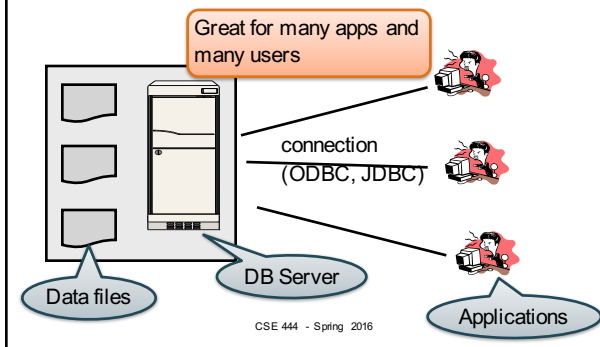
## DBMS Deployment: Local



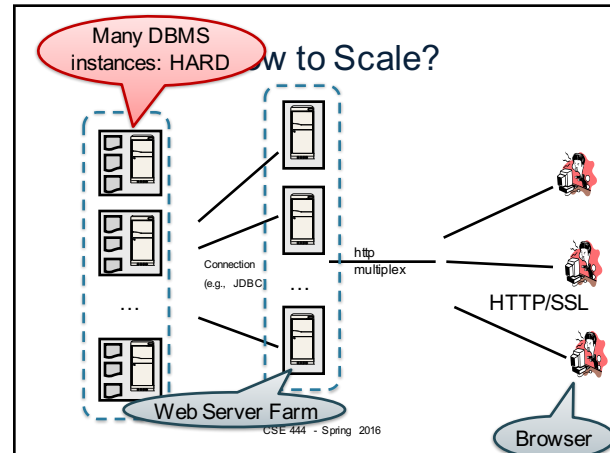
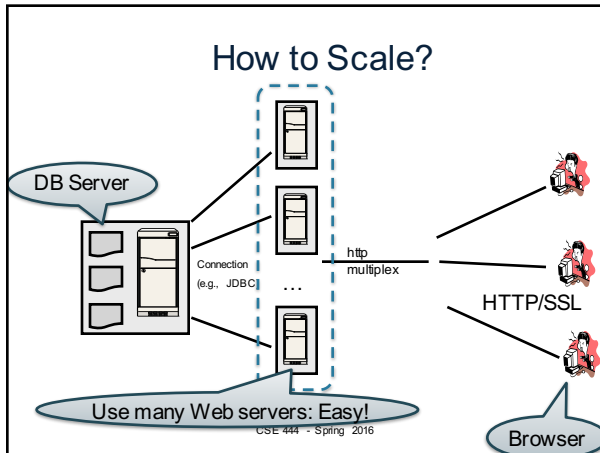
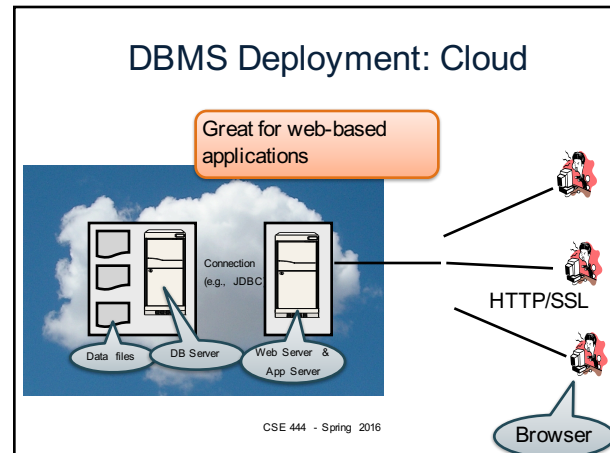
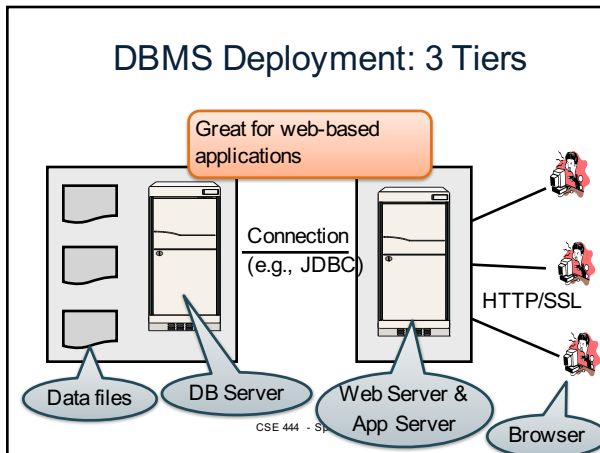
CSE 444 - Spring 2016

5

## DBMS Deployment: Client/Server



CSE 444 - Spring 2016

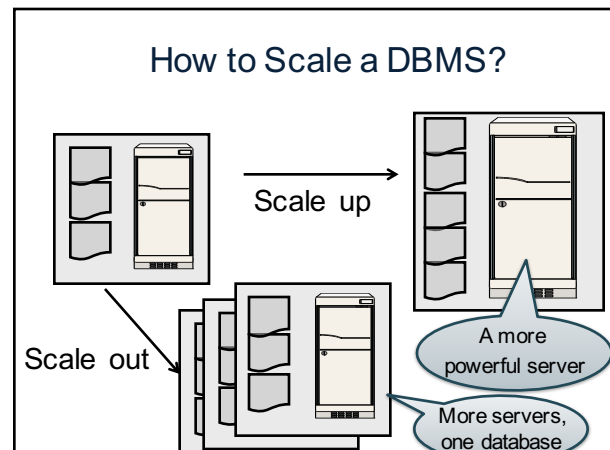


### How to Scale?

- 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

CSE 444 - Spring 2016

11



## What to scale?

- OLTP: Transactions per second
  - OLTP = Online Transaction Processing
- OLAP: Query response time
  - OLAP = Online Analytical Processing

CSE 444 - Spring 2016

13

## 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

CSE 444 - Spring 2016

14

## Scaling Single Query Response Time

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

CSE 444 - Spring 2016

15

## This Week: Focus on Scaling a Single Query

CSE 444 - Spring 2016

16

## Big Data

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

CSE 444 - Spring 2016

17

## 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

CSE 444 - Spring 2016

18

## 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

CSE 444 - Spring 2016

19

## 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
- **DATAAllegro** 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

21

## 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

CSE 444 - Spring 2016

22

## Performance Metrics for Parallel DBMSs

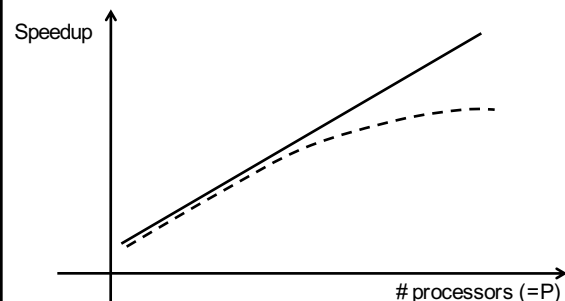
### Speedup

- 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")

CSE 444 - Spring 2016

23

## Linear v.s. Non-linear Speedup



CSE 444 - Spring 2016

24

## Performance Metrics for Parallel DBMSs

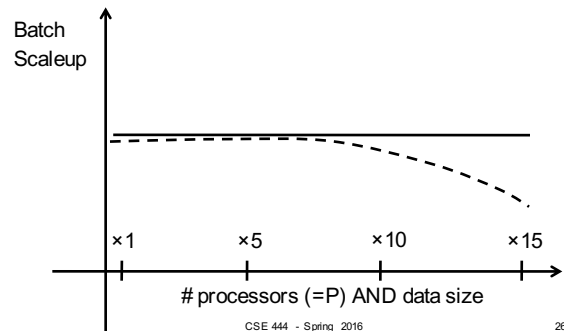
### Scaleup

- 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

CSE 444 - Spring 2016

25

## Linear v.s. Non-linear Scaleup



CSE 444 - Spring 2016

26

## Warning

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

CSE 444 - Spring 2016

27

## 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

CSE 444 - Spring 2016

28

## Three Architectures for Parallel DB

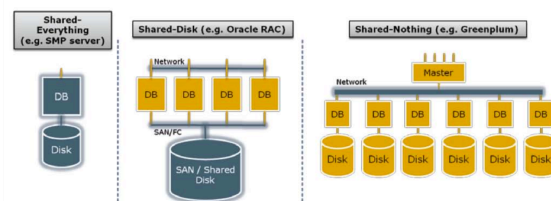
- Shared memory
- Shared disk
- Shared nothing

CSE 444 - Spring 2016

29

## Architectures for Parallel Databases

Figure 1 - Types of database architecture



From: Greenplum Database Whitepaper

SAN = "Storage Area Network"

CSE 444 - Spring 2016

30

## 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

CSE 444 - Spring 2016

31

## 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

CSE 444 - Spring 2016

32

## 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

33

## In Class

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

CSE 444 - Spring 2016

34

## Taxonomy for Parallel Query Evaluation

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

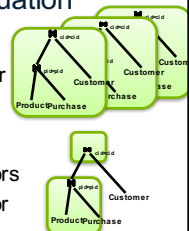


CSE 444 - Spring 2016

35

## 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



CSE 444 - Spring 2016

36

# 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

The diagram illustrates three types of parallel query evaluation using a query tree example. The query tree has a root node 'Join' with two children: 'ProductPurchase' and 'Customer'. The 'ProductPurchase' node has a child 'ProductPurchase' (labeled 'ProductPurchase' in the diagram), and the 'Customer' node has a child 'Customer' (labeled 'Customer' in the diagram). The 'ProductPurchase' node has a child 'ProductPurchase' (labeled 'ProductPurchase' in the diagram).

- Inter-query parallelism:** The query tree is replicated across multiple processors. Each processor runs a separate instance of the query.
- Inter-operator parallelism:** The query tree is replicated across multiple processors. Each processor runs a separate instance of the query.
- Intra-operator parallelism:** The query tree is replicated across multiple processors. Each processor runs a separate instance of the query.

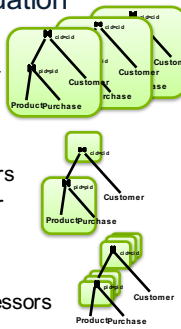
CSE 444 - Spring 2016

37

- **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

CSE 444 - Spring 2016

37



# 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

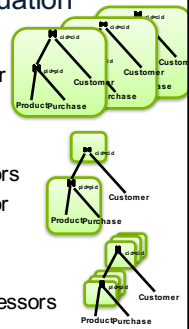
The diagram illustrates three levels of parallel query evaluation:

- Inter-query parallelism:** Multiple separate query trees, each in its own green box. Each tree has a root node 'q' and children 'ProductPurchase' and 'Customer'.
- Inter-operator parallelism:** A single query tree with some operators (like 'ProductPurchase') replicated in multiple green boxes. The root node 'q' has children 'ProductPurchase' and 'Customer'.
- Intra-operator parallelism:** A single query tree where a specific operator (like 'ProductPurchase') is further divided into multiple parallel sub-tasks within its box. The root node 'q' has children 'ProductPurchase' and 'Customer'.

We study only intra-operator parallelism: most scalable

- **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

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:**  $\gamma_{A \text{ sum}(B)}(R)$
- **Join:**  $R \bowtie S$

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

CSE 444 - Spring 2016 39

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

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

CSE 444 - Spring 2016

39

[illegible]

Servers:

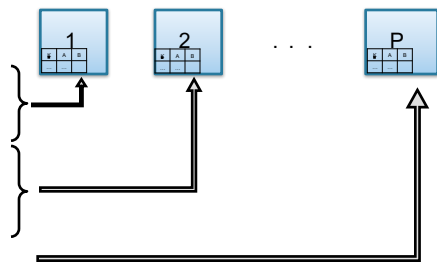
[illegible]

CSE 444 - Spring 2016

40

[illegible]

Servers:

[illegible]

---

CSE 444 - Spring 2016

41

# Horizontal Data Partitioning

Data:

K	A	B
...	...	

Servers:

1  

1	2	3
4	5	6

2  

1	2	3
4	5	6

...

P  

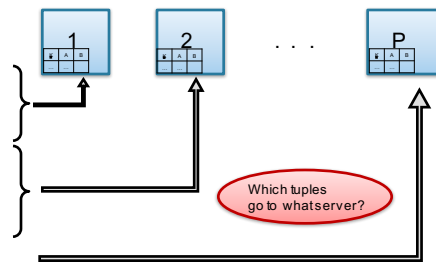
1	2	3
4	5	6

Which tuples go to what server?

CSE 444 - Spring, 2016

42

Servers:

[illegible]

---

CSE 444 - Spring 2016

42

## Horizontal Data Partitioning

- Relation  $R$  split into  $P$  chunks  $R_0, \dots, R_{P-1}$ , 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  $h(t.A) \bmod P$
- **Range based partitioning on attribute  $A$ :**
  - Tuple  $t$  to chunk  $i$  if  $v_{i-1} < t.A \leq v_i$

CSE 444 - Spring 2016

43

## Uniform Data v.s. Skewed Data

- 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  $A$

CSE 444 - Spring 2016

44

## Uniform Data v.s. Skewed Data

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

- **Block partition**

Uniform

- **Hash-partition**
  - On the key  $K$
  - On the attribute  $A$

Uniform Assuming uniform hash function

- **Range-partition**
  - On the key  $K$
  - On the attribute  $A$

CSE 444 - Spring 2016

45

## Uniform Data v.s. Skewed Data

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

- **Block partition**

Uniform

- **Hash-partition**
  - On the key  $K$
  - On the attribute  $A$

Uniform Assuming uniform hash function

- **Range-partition**
  - On the key  $K$
  - On the attribute  $A$

CSE 444 - Spring 2016

46

May be skewed  
E.g. when all records have the same value of the attribute  $A$ , then all records end up in the same partition

## Uniform Data v.s. Skewed Data

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

- **Block partition**

Uniform

- **Hash-partition**
  - On the key  $K$
  - On the attribute  $A$

Uniform Assuming uniform hash function

- **Range-partition**
  - On the key  $K$
  - On the attribute  $A$

CSE 444 - Spring 2016

47

May be skewed  
E.g. when all records have the same value of the attribute  $A$ , then all records end up in the same partition

May be skewed  
Difficult to partition the range of  $A$  uniformly.

## 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

CSE 444 - Spring 2016

48



## Horizontal Data Partitioning

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

CSE 444 - Spring 2016

49

## 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

CSE 444 - Spring 2016

50

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

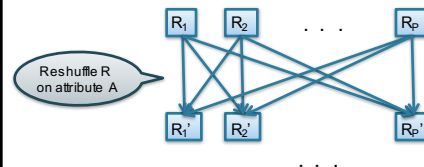
CSE 444 - Spring 2016

51

## Basic Parallel GroupBy

Data:  $R(K, A, B, C)$  -- hash-partitioned on  $K$

Query:  $\gamma_{A, \text{sum}(B)}(R)$



CSE 444 - Spring 2016

52

## Basic Parallel GroupBy

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

CSE 444 - Spring 2016

53

## Basic Parallel GroupBy

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

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

CSE 444 - Spring 2016

54

## Basic Parallel GroupBy

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

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

CSE 444 - Spring 2016

55

## Basic Parallel GroupBy

Can we do better?

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

CSE 444 - Spring 2016

56

## Basic Parallel GroupBy

Can we do better?

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

Distributive	Algebraic	Holistic
$\text{sum}(a_1 + a_2 + \dots + a_n) =$ $\text{sum}(\text{sum}(a_1 + a_2 + a_3) +$ $\text{sum}(a_4 + a_5 + a_6) +$ $\text{sum}(a_7 + a_8 + a_9))$	$\text{avg}(B) =$ $\text{sum}(B) / \text{count}(B)$	$\text{median}(B)$

CSE 444 - Spring 2016

57

## Parallel Join: $R \bowtie_{A=B} S$

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

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



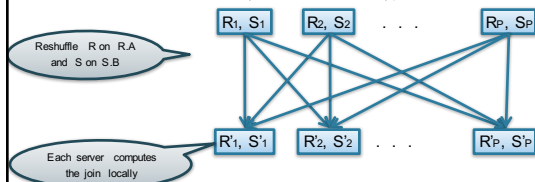
CSE 444 - Spring 2016

58

## Parallel Join: $R \bowtie_{A=B} S$

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

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



CSE 444 - Spring 2016

59

## Parallel Join: $R \bowtie_{A=B} S$

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

CSE 444 - Spring 2016

60

## Parallel Join: $R \bowtie_{A=B} S$

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 ?

CSE 444 - Spring 2016

61

## Parallel Join: $R \bowtie_{A=B} S$

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 ?
- **A:**  $(B(R)+B(S)) / P$

CSE 444 - Spring 2016

62

## Speedup and Scaleup

- Consider:
  - Query:  $\gamma_{Asum(C)}(R)$
  - 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?

CSE 444 - Spring 2016

63

## Speedup and Scaleup

- Consider:
  - Query:  $\gamma_{Asum(C)}(R)$
  - 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  $\frac{1}{2}$  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)

CSE 444 - Spring 2016

64

## Optimization for Small Relations

When joining  $R$  and  $S$

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

CSE 444 - Spring 2016

65

## 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

CSE 444 - Spring 2016

66

## 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} → [1,2] and [3,6]
- Eq-depth v.s. eq-width histograms

CSE 444 - Spring 2016

67

## Some Skew Handling Techniques

Create more partitions than nodes

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

CSE 444 - Spring 2016

68

## Some Skew Handling Techniques

Use subset-replicate (a.k.a. “skewedJoin”)

- Given  $R \bowtie_{A=B} S$
- 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

CSE 444 - Spring 2016

69

## 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

CSE 444 - Spring 2016

70

## Parallel Query Evaluation

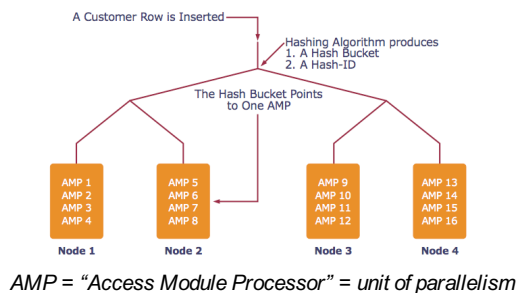
New operator: **Shuffle**

- Origin: **Exchange** operator from Volcano system
- Serves to re-shuffle data between processes
  - Handles data routing, buffering, and flow control
- Two parts: **ShuffleProducer** and **ShuffleConsumer**
- Producer:
  - Pulls data from child 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

CSE 444 - Spring 2016

71

## Example: Teradata – Loading



CSE 444 - Spring 2016

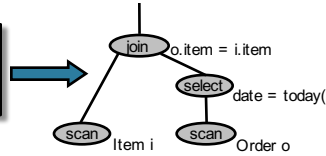
72

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

## Example: Teradata – Query Execution

Find all orders from today, along with the items ordered

```
SELECT *
FROM Order o, Line i
WHERE o.item = i.item
AND o.date = today()
```

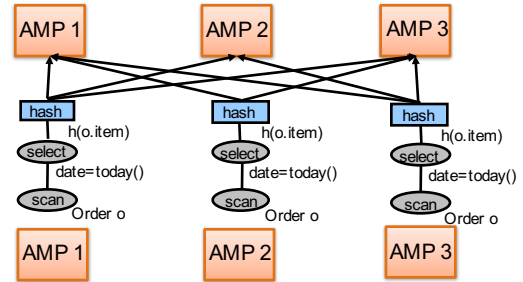
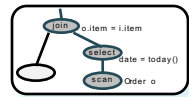


CSE 444 - Spring 2016

73

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

## Query Execution

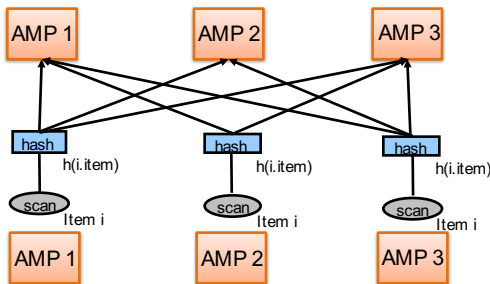
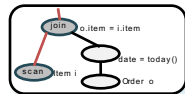


CSE 444 - Spring 2016

74

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

## Query Execution

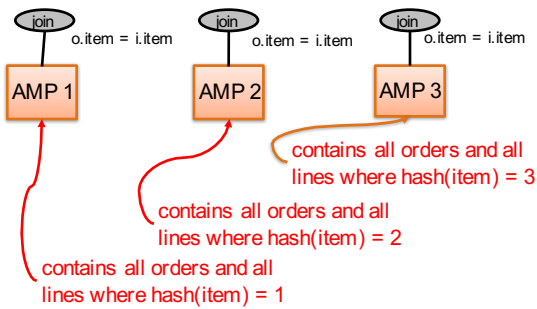


CSE 444 - Spring 2016

75

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

## Query Execution



CSE 444 - Spring 2016

76