# CompSci 516
# Database Systems

## Lecture 23-24
## Parallel DBMS
## Distributed DBMS
## NOSQL

Instructor: Sudeepa Roy

# Announcements (Thurs, 11/19)

- HW3/Mongo Due today!
- Keep working on projects!

# Reading Material

- [RG]
  - Parallel DBMS: Chapter 22.1-22.5
  - Distributed DBMS: Chapter 22.6 – 22.14

- [GUW]
  - Parallel DBMS and map-reduce: Chapter 20.1-20.2
  - Distributed DBMS: Chapter 20.3, 20.4.1-20.4.2, 20.5-20.6

- Other recommended readings:
  - Chapter 2 (Sections 1,2,3) of Mining of Massive Datasets, by Rajaraman and Ullman: http://i.stanford.edu/~ullman/mmds.html
  - Original Google MR paper by Jeff Dean and Sanjay Ghemawat, OSDI' 04: http://research.google.com/archive/mapreduce.html
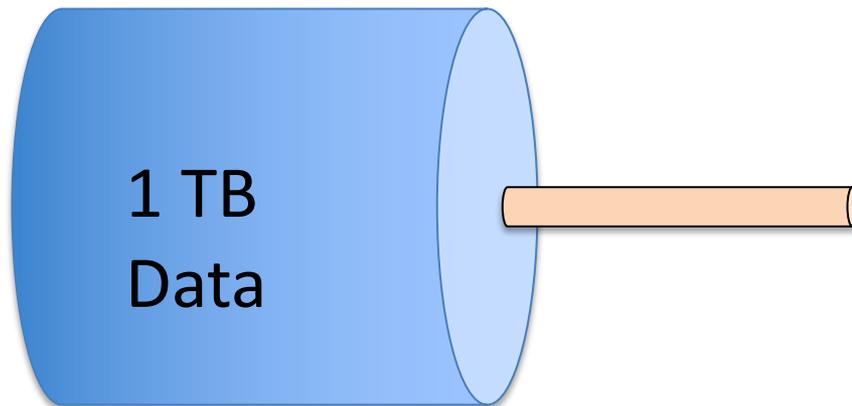
Acknowledgement:
The following slides have been created adapting the instructor material of the [RG] book provided by the authors Dr. Ramakrishnan and  Dr. Gehrke.
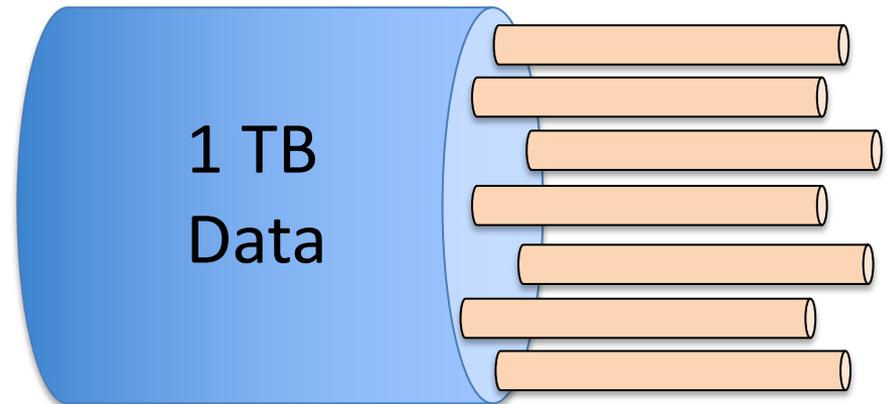
# Parallel DBMS

# Why Parallel Access To Data?

At 10 MB/s
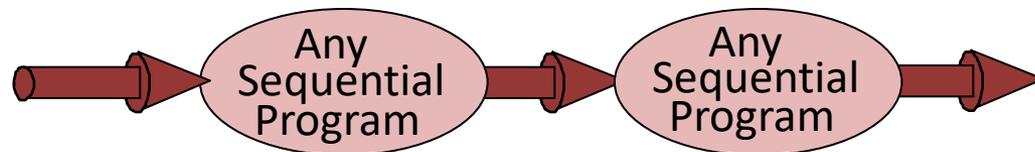1.2 days to scan

1,000 x parallel
1.5 minute to scan.

1 TB
Data

1 TB
Data

Parallelism:
divide a big problem
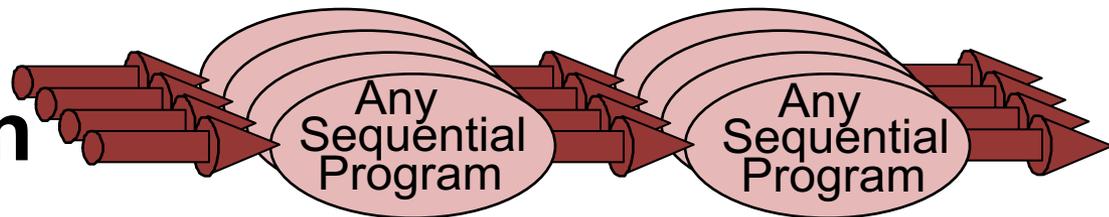into many smaller ones
to be solved in parallel.

# Parallel DBMS

- Parallelism is natural to DBMS processing
  - Pipeline parallelism: many machines each doing one step in a multi-step process.
  - Data-partitioned parallelism: many machines doing the same thing to different pieces of data.
  - Both are natural in DBMS!

**Pipeline**

**Partition**

**outputs split N ways, inputs merge M ways**

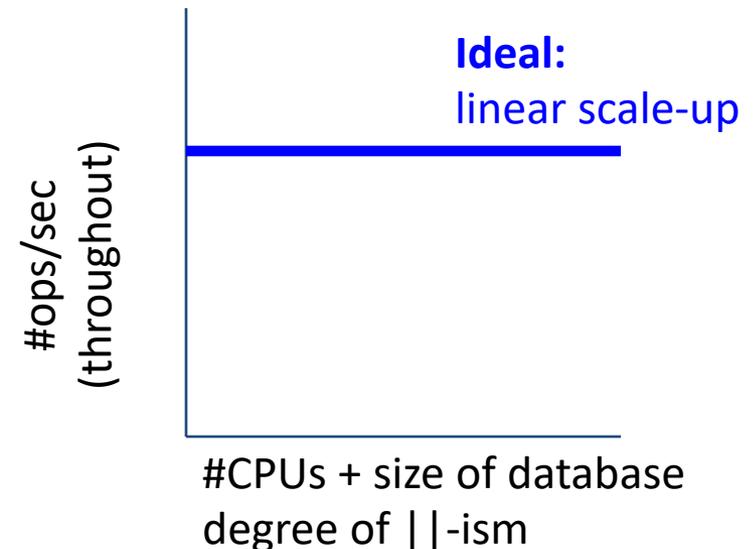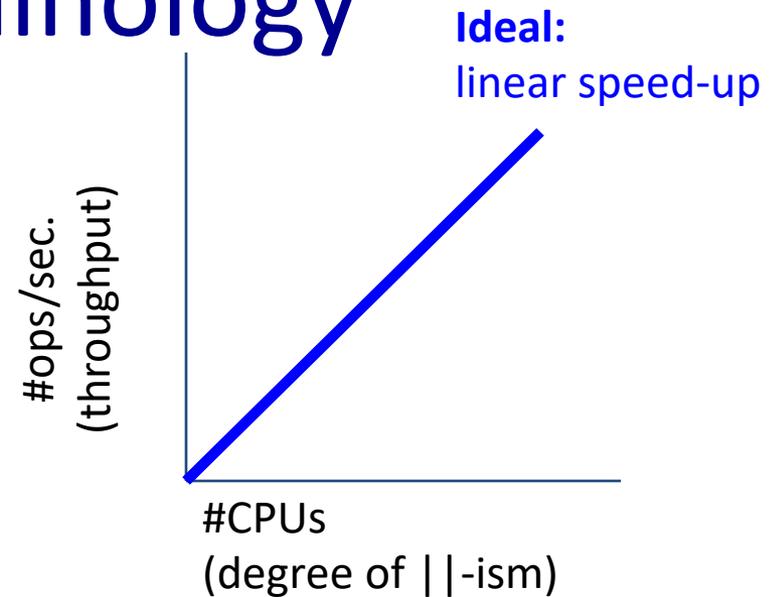# DBMS: The parallel Success Story

- 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

# Some || Terminology

**Ideal graphs**

- ## Speed-Up
  - – More resources means proportionally less time for given amount of data.

#ops/sec. (throughput)

#CPUs
(degree of ||-ism)

**Ideal:**
linear scale-up

- ## Scale-Up
  - – If resources increased in proportion to increase in data size, time is constant.

#ops/sec
(throughout)

#CPUs + size of database
degree of ||-ism

# Some || Terminology

## In practice

- Due to overhead in parallel processing

- Start-up cost

Starting the operation on many processor, might need to distribute data

- Interference

Different processors may compete for the same resources

- Skew

The slowest processor (e.g. with a huge fraction of data) may become the bottleneck

**Ideal:** linear speed-up

#ops/sec. (throughput)

**Actual:** sub-linear speed-up

#CPUs
(degree of ||-ism)

**Ideal:** linear scale-up

#ops/sec (throughput)

**Actual:** sub-linear scale-up

#CPUs + size of database
degree of ||-ism

# Basics of Parallelism

- Units: a collection of processors
  - assume always have local cache
  - may or may not have local memory or disk (next)

- A communication facility to pass information among processors
  - a shared bus or a switch

- Different architecture
  - Whether memory AND/OR disk are shared

# Shared Memory

- Easy to program
- Expensive to build
- Low communication overhead: shared mem.
- Difficult to scaleup (memory contention)

P    P    P

## Interconnection Network

shared memory

## Global Shared Memory

D    D    D

# Shared Disk

- Trade-off but still interference like shared-memory (contention of memory and nw bandwidth)

local memory

| M | M | | M |
|---|---|---|---|

| P | P | | P |

## Interconnection Network

shared disk

| D | D | | D |

# Shared Nothing

- Hard to program and design parallel algos
- Cheap to build
- Easy to scaleup and speedup
- Considered to be the best architecture
- We will assume this architecture!

local memory and disk

no two CPU can access the same storage area

all communication through a network connection

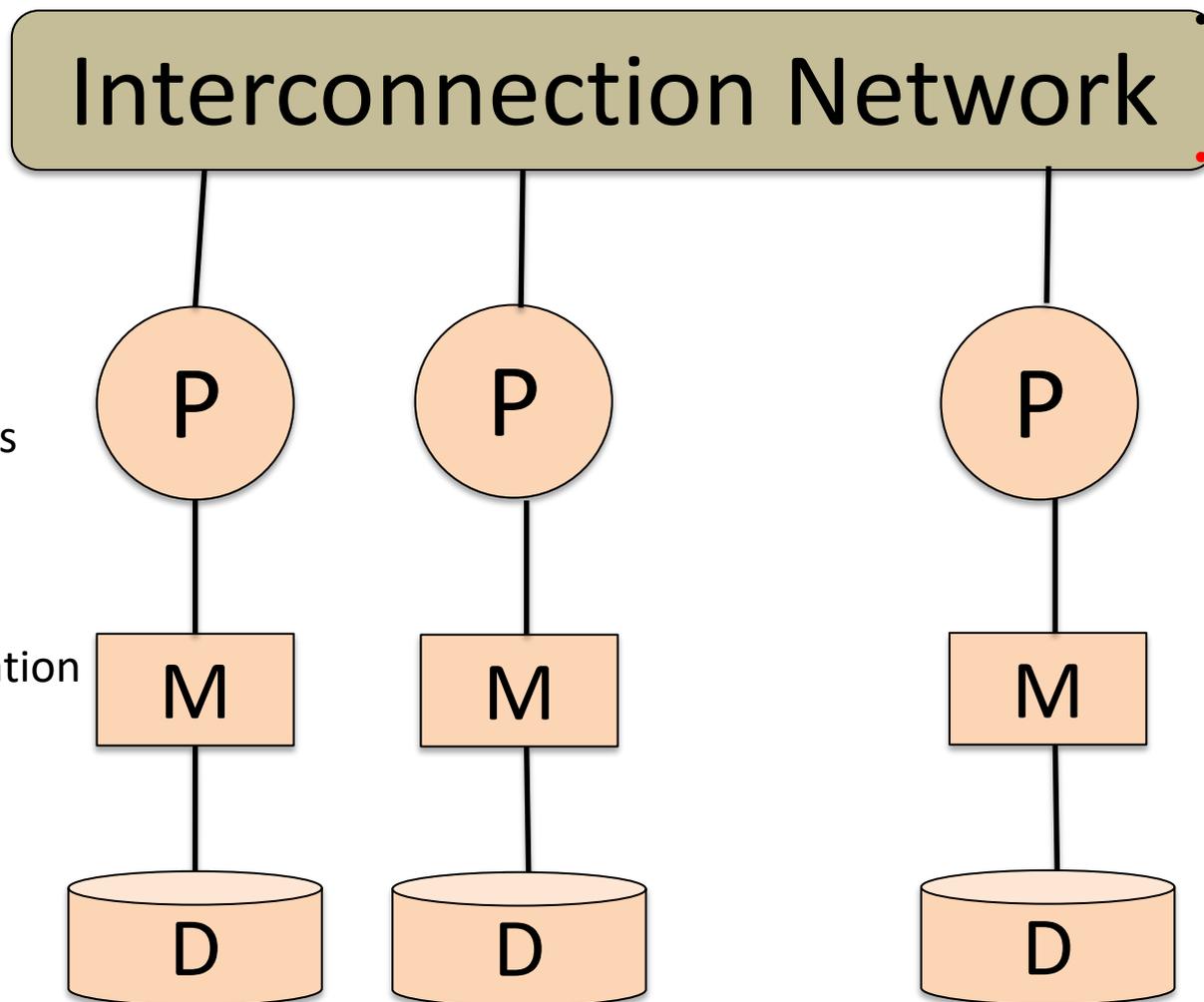## Interconnection Network

P    P    P

M    M    M

D    D    D

# Different Types of DBMS Parallelism

- ## Intra-operator parallelism
  - get all machines working to compute a given operation (scan, sort, join)
  - OLAP (decision support)

- ## Inter-operator parallelism
  - each operator may run concurrently on a different site (exploits pipelining)
  - For both OLAP and OLTP

- ## Inter-query parallelism
  - different queries run on different sites
  - For OLTP

- ## We'll focus on intra-operator parallelism

Ack:
Slide by Prof. Dan Suciu

# Data Partitioning

**Horizontally Partitioning a table** (why horizontal?):

| Range-partition | Hash-partition | Block-partition or Round Robin |
|---|---|---|



- Good for equijoins, range queries, group-by
- Can lead to data skew

- Good for equijoins
- But only if hashed on that attribute
- Can lead to data skew

- Send i-th tuple to i-mod-n processor
- Good to spread load
- Good when the entire relation is accessed

Shared disk and memory less sensitive to partitioning, Shared nothing benefits from "good" partitioning

# Best serial plan may not be best ||

- Why?

- Trivial counter-example:
  - Table partitioned with local secondary index at two nodes
  - Range query: all of node 1 and 1% of node 2.
  - Node 1 should do a scan of its partition.
  - Node 2 should use secondary index.

**Table Scan**

**Index Scan**

A..M

N..Z

# 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 <u>block partitioned </u>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**

R(a, b)

SELECT a, max(b) as topb
FROM R
WHERE a > 0
GROUP BY a

| Machine 1 | Machine 2 | Machine 3 |
|---|---|---|
| 1/3 of R | 1/3 of R | 1/3 of R |

R(a, b)

If more than one relation on a machine, then "scan S", "scan R" etc

```
  scan              scan              scan

Machine 1         Machine 2         Machine 3

1/3 of R          1/3 of R          1/3 of R
```

R(a, b)

SELECT a, max(b) as topb
FROM R
WHERE a > 0
GROUP BY a

$\sigma_{a>0}$

scan

Machine 1

1/3 of R

$\sigma_{a>0}$

scan

Machine 2

1/3 of R

$\sigma_{a>0}$

scan

Machine 3

1/3 of R

R(a, b)

SELECT a, max(b) as topb
FROM R
WHERE a > 0
GROUP BY a



$\gamma_{a, max(b) -> b}$

$\sigma_{a>0}$

scan

Machine 1

1/3 of R

$\gamma_{a, max(b) -> b}$

$\sigma_{a>0}$

scan

Machine 2

1/3 of R

$\gamma_{a, max(b) -> b}$

$\sigma_{a>0}$

scan

Machine 3

1/3 of R

R(a, b)

SELECT a, max(b) as topb
FROM R
WHERE a > 0
GROUP BY a



**Hash on a**

$\gamma_{a, \text{ max(b)-> b}}$

$\sigma_{a>0}$

**scan**

Machine 1

1/3 of R

**Hash on a**

$\gamma_{a, \text{ max(b)-> b}}$

$\sigma_{a>0}$

**scan**

Machine 2

1/3 of R

**Hash on a**

$\gamma_{a, \text{ max(b)-> b}}$

$\sigma_{a>0}$

**scan**

Machine 3

1/3 of R
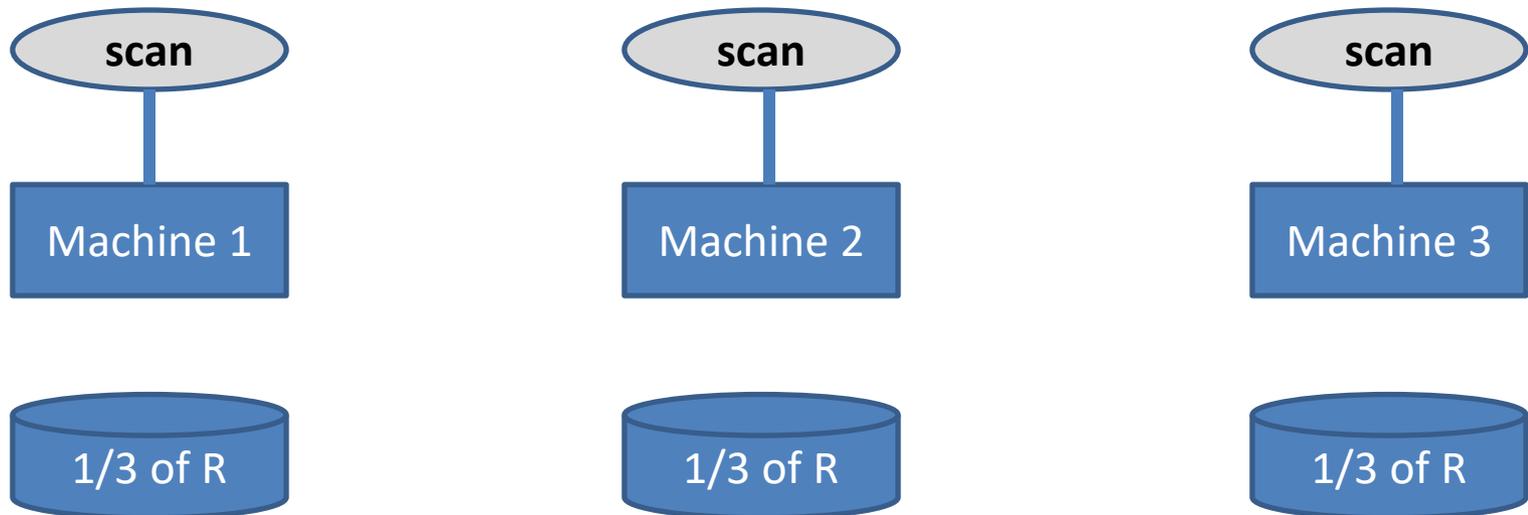
R(a, b)

SELECT a, max(b) as topb     FROM R
WHERE a > 0               GROUP BY a

**Hash on a**     **Hash on a**     **Hash on a**

$\gamma_{a,\ max(b)\ ->\ b}$     $\gamma_{a,\ max(b)\ ->\ b}$     $\gamma_{a,\ max(b)\ ->\ b}$

$\sigma_{a>0}$     $\sigma_{a>0}$     $\sigma_{a>0}$

**scan**     **scan**     **scan**

Machine 1     Machine 2     Machine 3

1/3 of R     1/3 of R     1/3 of R

R(a, b)

$\gamma_{a, \, max(b)->topb}$

$\gamma_{a, \, max(b)->topb}$

$\gamma_{a, \, max(b)->topb}$

Hash on a

Hash on a

Hash on a

$\gamma_{a, \, max(b)-> b}$

$\gamma_{a, \, max(b)-> b}$

$\gamma_{a, \, max(b)-> b}$

$\sigma_{a>0}$

$\sigma_{a>0}$

$\sigma_{a>0}$

scan

scan

scan

Machine 1

Machine 2

Machine 3

1/3 of R

1/3 of R

1/3 of R

# Benefit of hash-partitioning

- What would change if we hash-partitioned R on R.a before executing the same query on the previous parallel DBMS and MR

**Prev: block-partition**

SELECT a, max(b) as topb    FROM R
WHERE a > 0            GROUP BY a

$\gamma_{a, max(b)->topb}$    $\gamma_{a, max(b)->topb}$    $\gamma_{a, max(b)->topb}$

Hash on a        Hash on a        Hash on a

$\gamma_{a, max(b)-> b}$    $\gamma_{a, max(b)-> b}$    $\gamma_{a, max(b)-> b}$

$\sigma_{a>0}$        $\sigma_{a>0}$        $\sigma_{a>0}$

scan            scan            scan

Machine 1        Machine 2        Machine 3

1/3 of R        1/3 of R        1/3 of R

SELECT a, max(b) as topb
FROM R
WHERE a > 0
GROUP BY a

- It would avoid the data re-shuffling phase
- It would compute the aggregates locally

**Hash-partition on a for R(a, b)**
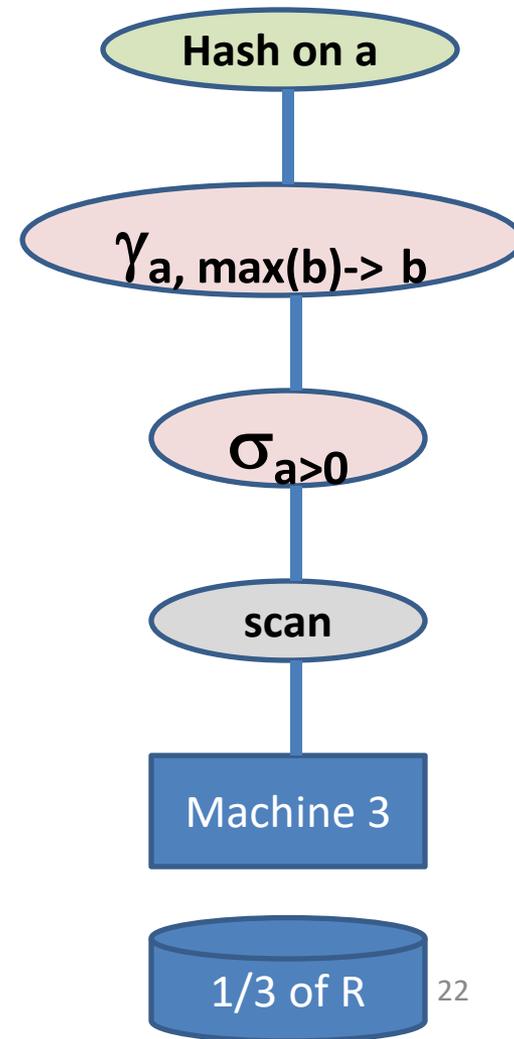
SELECT a, max(b) as topb     FROM R
WHERE a > 0               GROUP BY a

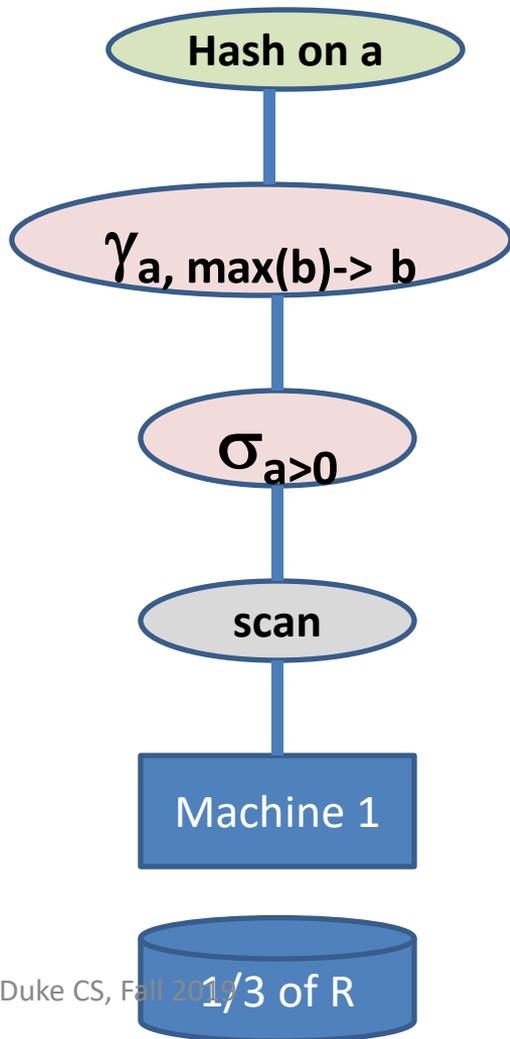$\gamma_{\text{a, max(b)->topb}}$      $\gamma_{\text{a, max(b)->topb}}$      $\gamma_{\text{a, max(b)->topb}}$

$\sigma_{a>0}$          $\sigma_{a>0}$          $\sigma_{a>0}$

scan          scan          scan

Machine 1          Machine 2          Machine 3

1/3 of R          1/3 of R          1/3 of R

# Distributed DBMS

# Parallel vs. Distributed DBMS

## Parallel DBMS

- Parallelization of various operations
  - e.g. loading data, building indexes, evaluating queries

- Data may or may not be distributed initially

- Distribution is governed by performance consideration

## Distributed DBMS

- Data is physically stored across different sites
  - Each site is typically managed by an independent DBMS

- Location of data and autonomy of sites have an impact on Query opt., Conc. Control and recovery

- Also governed by other factors:
  - increased availability for system crash
  - local ownership and access

# Topics in Distributed DBMS

- Architecture

- Data Storage

- Query Execution

- Transactions – updates

- Recovery – Two Phase Commit (2PC)

- A brief overview / examples of all these

# Distributed Data Independence

- Users should not have to know where data is located
  - no need to know the locations of references relations, their copies or fragments (later)
  - extends Physical and Logical Data Independence principles

- Queries spanning multiple sites should be optimized in a cost-based manner
  - taking into account communication costs and differences in local computation costs

# Distributed DBMS Architectures

- Three alternative approaches

1. Client-Server
   – Client: user interace, server: executes queries

2. Collaborating Server
   – All are of the same status

3. Middleware
   – Good for integrating legacy systems, middleware coordinates, individual server executes local queries

# Storing Data in a Distributed DBMS

- A single relation may be partitioned or fragmented across several sites
  - typically at sites where they are most often accessed

- The data can be replicated as well
  - when the relation is in high demand or for robustness

- Horizontal:
  - Usually disjoint
  - Can often be identified by a selection query
    - employees in a city – locality of reference
  - To retrieve the full relation, need a union

- Vertical:
  - Identified by projection queries
  - Typically unique TIDs added to each tuple
  - TIDs replicated in each fragments
  - Ensures that we have a Lossless Join

**TID**

| TID | | | | | |
|-----|--|--|--|--|--|
| t1  | | | | | |
| t2  | | | | | |
| t3  | | | | | |
| t4  | | | | | |

# Joins in a Distributed DBMS

- Can be very expensive if relations are stored at different sites

1. Fetch as needed

Sailors as outer – for each S page, fetch all R pages from Paris
if cached at London, each R page fetched once

2. Ship to one site

Ship Sailor to Paris

Unnecessary shipping
Not all tuples used

3. Semi-join

4. Bloom join

**LONDON**

**Sailors (S)**

500 pages

**PARIS**

**Reserves (R)**

1000 pages

# Semijoin

- Suppose want to ship R to London and then do join with S at London. Instead,

1. At London, project S onto join columns and ship this to Paris
   - Here foreign keys, but could be arbitrary join

2. At Paris, join S-projection with R
   - Result is called reduction of Reserves w.r.t. Sailors (only these tuples are needed)

3. Ship reduction of R to back to London

4. At London, join S with reduction of R

- Tradeoff the cost of computing and shipping projection for cost of shipping full R relation
  - Especially useful if there is a selection on Sailors, and answer desired at London

# Bloomjoin

**LONDON**

Sailors (S)

500 pages

**PARIS**

Reserves (R)

1000 pages

- Similar idea like semi-join
- Suppose want to ship R to London and then do join with S at London (like semijoin)

1. At London, compute a bit-vector of some size k:
   - Hash column values into range 0 to k-1
   - If some tuple hashes to p, set bit p to 1 (p from 0 to k-1)
   - Ship bit-vector to Paris

2. At Paris, hash each tuple of R similarly
   - discard tuples that hash to 0 in S's bit-vector
   - Result is called reduction of R w.r.t S

3. Ship "bit-vector-reduced" R to London

4. At London, join S with reduced R

- Bit-vector cheaper to ship, almost as effective
   - the size of the reduction of R shipped back can be larger. Why?

# Distributed Query Optimization

- Similar to centralized optimization, but have differences
    1. Communication costs must be considered
    2. Local site autonomy must be respected
    3. New distributed join methods should be considered

- Query site constructs global plan, with suggested local plans describing processing at each site
    - If a site can improve suggested local plan, free to do so

# Updating Distributed Data

- Synchronous Replication: All copies of a modified relation (or fragment) must be updated before the modifying transaction commits
  - Always updated but expensive commit protocols (2PC – soon!)
  - By "voting" - e.g., 10 copies; 7 written for update; 4 copies read (why 4?)
  - Read-any Write-all (special case of voting, why not write-any read all?)

- Asynchronous Replication:  Copies of a modified relation are only periodically updated; different copies may get out-of-sync in the meantime
  - More efficient – many current products follow this approach
  - Primary site (one master copy) or peer-to-peer (multiple master copies)

# Distributed Locking

- How do we manage locks for objects across many sites?

1. Centralized: One site does all locking
   - Vulnerable to single site failure

2. Primary Copy: All locking for an object done at the primary copy site
   - Reading requires access to locking site as well as site where the object copy is stored

3. Fully Distributed: Locking for a copy done at site where the copy is stored
   - Locks at all sites while writing an object (unlike previous two)
   - May lead to "undetected" or "missing" "global deadlock" due to delay in information propagation
   - Timeout or hierarchical detection
     - e.g. sites (every 10 sec)-> sites in a state (every min)-> sites in a country (every 10 min) -> global waits for graph. Intuition: more deadlocks are likely across closely related sites

# Distributed Recovery

- Two new issues:
  - New kinds of failure, e.g., links and remote sites
  - If "sub-transactions" of a transaction execute at different sites, all or none must commit
  - Need a commit protocol to achieve this
  - Most widely used: Two Phase Commit (2PC)

- A log is maintained at each site
  - as in a centralized DBMS
  - commit protocol actions are additionally logged

# Two-Phase Commit (2PC)

- Site at which transaction originates is coordinator

- Other sites at which it executes are subordinates

    - w.r.t. coordination of this transaction

Example on whiteboard

# When a transaction wants to commit – 1/5

1. Coordinator sends prepare message to each subordinate

# When a transaction wants to commit – 2/5

2. Subordinate receives the prepare message
    a) decides whether to abort or commit its subtransaction
    b) force-writes an <span style="color:red">abort</span> or <span style="color:red">prepare</span> log record
    c) then sends a <span style="color:red">no</span> or <span style="color:red">yes</span> message to coordinator

# When a transaction wants to commit – 3/5

3. If coordinator gets unanimous yes votes from all subordinates
   a) it force-writes a commit log record
   b) then sends commit message to all subs

Else (if receives a no message or no response from some subordinate),
   a) it force-writes abort log record
   b) then sends abort messages

4.  Subordinates force-write abort/commit log record based on message they get

   a)  then send ack message to coordinator

   b)  If commit received, commit the subtransaction

   c)  write an end record

# When a transaction wants to commit – 5/5

5.  After the coordinator receives ack from all subordinates,
    – writes end log record

Transaction is officially committed when the coordinator's commit log record reaches the disk
    – subsequent failures cannot affect the outcomes

# Comments on 2PC

- Two rounds of communication
  - first, voting
  - then, termination
  - Both initiated by coordinator
- Any site (coordinator or subordinate) can unilaterally decide to abort a transaction
  - but unanimity/consensus needed to commit
- Every message reflects a decision by the sender
  - to ensure that this decision survives failures, it is first recorded in the local log and is force-written to disk
- All commit protocol log records for a transaction contain tid and Coordinator-id
  - The coordinator's abort/commit record also includes ids of all subordinates.

# Restart After a Failure at a Site – 1/4

- Recovery process is invoked after a sites comes back up after a crash
  - reads the log and executes the commit protocol
  - the coordinator or a subordinate may have a crash
  - one site can be the coordinator some transaction and subordinates for others

# Restart After a Failure at a Site – 2/4

- If we have a commit or abort log record for transaction T, but not an end record, must redo/undo T respectively

  - If this site is the coordinator for T (from the log record), keep sending commit/abort messages to subs until acks received

  - then write an end log record for T

# Restart After a Failure at a Site – 3/4

- If we have a <span style="color:blue">prepare</span> log record for transaction T, but not <span style="color:red">commit/abort</span>
    - This site is a subordinate for T
    - Repeatedly contact the coordinator to find status of T
    - Then write commit/abort log record
    - Redo/undo T
    - and write end log record

# Restart After a Failure at a Site – 4/4

- If we don't have even a prepare log record for T
  - T was not voted to commit before crash
  - unilaterally abort and undo T
  - write an end record
- No way to determine if this site is the coordinator or subordinate
  - If this site is the coordinator, it might have sent prepare messages
  - then, subs may send yes/no message – coordinator is detected – ask subordinates to abort

# Blocking

- If coordinator for transaction T fails, subordinates who have voted yes cannot decide whether to commit or abort T until coordinator recovers.
  - T is blocked
  - Even if all subordinates know each other (extra overhead in prepare message) they are blocked unless one of them voted no
- Note: even if all subs vote yes, the coordinator then can give a no vote, and decide later to abort!

# Link and Remote Site Failures

- If a remote site does not respond during the commit protocol for transaction T, either because the site failed or the link failed:

  - If the current site is the coordinator for T, should abort T

  - If the current site is a subordinate, and has not yet voted yes, it should abort T

  - If the current site is a subordinate and has voted yes, it is blocked until the coordinator responds

  - needs to periodically contact the coordinator until receives a reply

# Observations on 2PC

- Ack messages used to let coordinator know when it can "forget" a transaction; until it receives all acks, it must keep T in the transaction Table

- If coordinator fails after sending prepare messages but before writing commit/abort log records, when it recovers, it aborts the transaction

- If a subtransaction does no updates, its commit or abort status is irrelevant

# NoSQL



- Optional reading:
    - Cattell's paper (2010-11)
    - Warning! some info will be outdated
    - see webpage http://cattell.net/datastores/ for updates and more pointers

CompSci 516: Database Systems

# NOSQL

- Many of the new systems are referred to as "NoSQL" data stores
  - MongoDB, CouchDB, VoltDB, Dynamo, Membase, ….
- NoSQL stands for "Not Only SQL" or "Not Relational"
  - not entirely agreed upon
- NoSQL = "new" database systems
  - not typically RDBMS
  - relax on some requirements, gain efficiency and scalability
- New systems choose to use/not use several concepts we learnt so far
  - You may find systems that use multi-version Concurrency Control (MVCC) or, asynchronous replication

| OLTP (Online Transaction Processing) | Data Warehousing/OLAP (On Line Analytical Processing) |
| --- | --- |
| Mostly updates | Mostly reads |
| Applications:<br>Order entry, sales update, banking transactions | Applications:<br>Decision support in industry/organization |
| Detailed, up-to-date data | Summarized, historical data<br>(from multiple operational db, grows over time) |
| Structured, repetitive, short tasks | Query intensive, ad hoc, complex queries |
| Each transaction reads/updates only a few tuples (tens of) | Each query can access many records, and perform many joins, scans, aggregates |
| MB-GB data | GB-TB data |
| Typically clerical users | Decision makers, analysts as users |
| Important:<br>Consistency, recoverability, Maximizing tr. throughput | Important:<br>Query throughput<br>Response times |

# Applications of New Systems

- Designed to scale simple "OLTP"-style application loads
  - to do updates as well as reads
  - in contrast to traditional DBMSs and data warehouses
  - to provide good horizontal scalability for simple read/write database operations distributed over many servers

- Originally motivated by Web 2.0 applications
  - these systems are designed to scale to thousands or millions of users

# NoSQL: Six Key Features

1.  the ability to horizontally scale "simple operations" throughput over many servers

2.  the ability to replicate and to distribute (partition) data over many servers

3.  a simple call level interface or protocol (in contrast to SQL binding)

4.  a weaker concurrency model than the ACID transactions of most relational (SQL) database systems

5.  efficient use of distributed indexes and RAM for data storage

6.  the ability to dynamically add new attributes to data records

# BASE (not ACID ☺)

- Recall ACID for RDBMS desired properties of transactions:
  - Atomicity, Consistency, Isolation, and Durability

- NOSQL systems typically do not provide ACID

- Basically Available

- Soft state

- Eventually consistent

# ACID vs. BASE

- The idea is that by giving up ACID constraints, one can achieve much higher performance and scalability

- The systems differ in how much they give up
  - e.g. most of the systems call themselves "eventually consistent", meaning that updates are eventually propagated to all nodes
  - but many of them provide mechanisms for some degree of consistency, such as multi-version concurrency control (MVCC)

# "CAP" "Theorem"

- Often Eric Brewer's CAP theorem cited for NoSQL

- A system can have only two out of three of the following properties:
  - Consistency,
  - Availability
  - Partition-tolerance

- The NoSQL systems generally give up consistency
  - However, the trade-offs are complex

# What is different in NOSQL systems

- When you study a new NOSQL system, notice how it differs from RDBMS in terms of

1. Concurrency Control

2. Data Storage Medium

3. Replication

4. Transactions

# Choices in NOSQL systems:
# 1. Concurrency Control

a) Locks
  – some systems provide one-user-at-a-time read or update locks
  – MongoDB provides locking at a field level

b) MVCC

c) None
  – do not provide atomicity
  – multiple users can edit in parallel
  – no guarantee which version you will read

d) ACID
  – pre-analyze transactions to avoid conflicts
  – no deadlocks and no waits on locks

# Choices in NOSQL systems:
# 2. Data Storage Medium

a)  Storage in RAM

– snapshots or replication to disk

– poor performance when overflows RAM

b)  Disk storage

– caching in RAM

# Choices in NOSQL systems: 3. Replication

- whether mirror copies are always in sync

a) Synchronous

b) Asynchronous

  – faster, but updates may be lost in a crash

c) Both

  – local copies synchronously, remote copies asynchronously

# Choices in NOSQL systems:
# 4. Transaction Mechanisms

a) support

b) do not support

c) in between

   – support local transactions only within a single object or "shard"

   – shard = a horizontal partition of data in a database

# Comparison from Cattell's paper (2011)

| System | Conc Contol | Data Storage | Repli-cation | Tx |
|---|---|---|---|---|
| Redis | Locks | RAM | Async | N |
| Scalaris | Locks | RAM | Sync | L |
| Tokyo | Locks | RAM or disk | Async | L |
| Voldemort | MVCC | RAM or BDB | Async | N |
| Riak | MVCC | Plug-in | Async | N |
| Membrain | Locks | Flash + Disk | Sync | L |
| Membase | Locks | Disk | Sync | L |
| Dynamo | MVCC | Plug-in | Async | N |
| SimpleDB | None | S3 | Async | N |
| MongoDB | Locks | Disk | Async | N |
| Couch DB | MVCC | Disk | Async | N |

| | | | | |
|---|---|---|---|---|
| Terrastore | Locks | RAM+ | Sync | L |
| HBase | Locks | Hadoop | Async | L |
| HyperTable | Locks | Files | Sync | L |
| Cassandra | MVCC | Disk | Async | L |
| BigTable | Locks+stamps | GFS | Sync+Async | L |
| PNUTs | MVCC | Disk | Async | L |
| MySQL Cluster | ACID | Disk | Sync | Y |
| VoltDB | ACID, no lock | RAM | Sync | Y |
| Clustrix | ACID, no lock | Disk | Sync | Y |
| ScaleDB | ACID | Disk | Sync | Y |
| ScaleBase | ACID | Disk | Async | Y |
| NimbusDB | ACID, no lock | Disk | Sync | Y |

# Data Store Categories

- The data stores are grouped according to their data model

- Key-value Stores:
  - store values and an index to find them based on a programmer- defined key
  - e.g. Project Voldemort, Riak, Redis, Scalaris, Tokyo Cabinet, Memcached/Membrain/Membase

- Document Stores:
  - store documents, which are indexed, with a simple query mechanism
  - e.g. Amazon SimpleDB, CouchDB, MongoDB, Terrastore

- Extensible Record Stores:
  - store extensible records that can be partitioned vertically and horizontally across nodes ("wide column stores")
  - e.g. Hbase, HyperTable, Cassandra, Yahoo's PNUTS

- "New" Relational Databases:
  - store (and index and query) tuples, e.g. the new RDBMSs that provide horizontal scaling
  - e.g. MySQL Cluster, VoltDB, Clustrix, ScaleDB, ScaleBase, NimbusDB, Google Megastore (a layer on BigTable)

# RDBMS benefits

- Relational DBMSs have taken and retained majority market share over other competitors in the past 30 years

- While no "one size fits all" in the SQL products themselves, there is a common interface with SQL, transactions, and relational schema that give advantages in training, continuity, and data interchange

- Successful relational DBMSs have been built to handle other specific application loads in the past:
  - read-only or read-mostly data warehousing, OLTP on multi-core multi-disk CPUs, in-memory databases, distributed databases, and now horizontally scaled databases

# NoSQL benefits

- We haven't yet seen good benchmarks showing that RDBMSs can achieve scaling comparable with NoSQL systems like Google's BigTable

- If you only require a lookup of objects based on a single key, then a key-value/document store may be adequate and probably easier to understand than a relational DBMS

- Some applications require a flexible schema

- A relational DBMS makes "expensive" (multi-node multi-table) operations "too easy"
  - NoSQL systems make them impossible or obviously expensive for programmers

- The new systems are slowly gaining market shares too
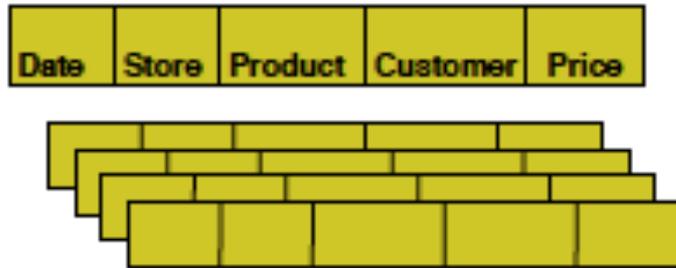
# Column Store

# Row vs. Column Store

- Row store
  - store all attributes of a tuple together
  - storage like "row-major order" in a matrix

- Column store
  - store all rows for an attribute (column) together
  - storage like "column-major order" in a matrix

- e.g.
  - MonetDB, Vertica (earlier, C-store), SAP/Sybase IQ, Google Bigtable (with column groups)
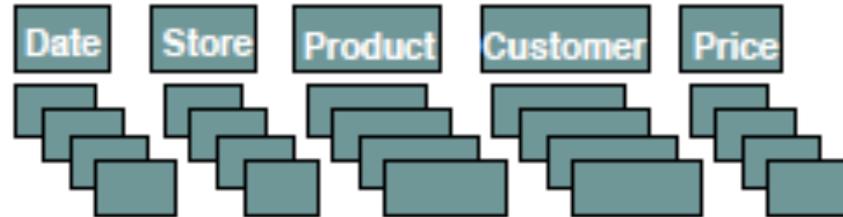
# What is a column-store?

## row-store

| Date | Store | Product | Customer | Price |
|------|-------|---------|----------|-------|

## column-store

| Date | | Store | | Product | | Customer | | Price |

+ easy to add/modify a record

- might read in unnecessary data

+ only need to read in relevant data

- tuple writes require multiple accesses

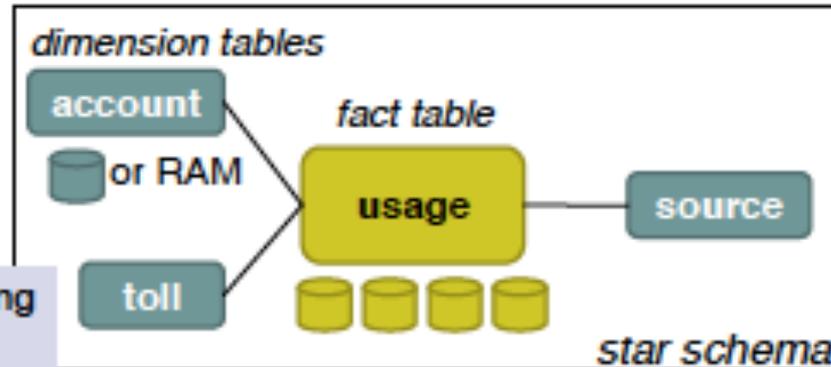=> *suitable for read-mostly, read-intensive, large data repositories*

Ack: Slide from VLDB 2009 tutorial on Column store

# Telco Data Warehousing example

1  **Typical DW installation**

1  **Real-world example**

"One Size Fits All? - Part 2: Benchmarking Results" Stonebraker et al. CIDR 2007

dimension tables

account
or RAM

fact table

usage

source

toll

star schema

```
QUERY 2
SELECT account.account_number,
sum (usage.toll_airtime),
sum (usage.toll_price)
FROM usage, toll, source, account
WHERE usage.toll_id = toll.toll_id
AND usage.source_id = source.source_id
AND usage.account_id = account.account_id
AND toll.type_ind in ('AE'. 'AA')
AND usage.toll_price > 0
AND source.type != 'CIBER'
AND toll.rating_method = 'IS'
AND usage.invoice_date = 20051013
GROUP BY account.account_number
```

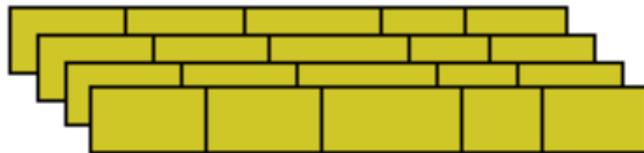|  | Column-store | Row-store |
|---|---|---|
| Query 1 | 2.06 | 300 |
| Query 2 | 2.20 | 300 |
| Query 3 | 0.09 | 300 |
| Query 4 | 5.24 | 300 |
| Query 5 | 2.88 | 300 |

**Why? Three main factors (next slides)**

Ack: Slide from  VLDB 2009 tutorial on Column store

# Telco example explained (1/3): *read efficiency*

## row store

read pages containing entire rows

one row = 212 columns!

is this typical? (it depends)

**What about vertical partitioning? (it does not work with ad-hoc queries)**

## column store

read only columns needed

in this example: 7 columns

caveats:
- "select * " not any faster
- clever disk prefetching
- clever tuple reconstruction

Ack: Slide from VLDB 2009 tutorial on Column store

# Telco example explained (2/3): compression efficiency

- Columns compress better than rows
  - Typical row-store compression ratio  1 : 3
  - Column-store 1 : 10

- Why?
  - Rows contain values from different domains
    => more entropy, difficult to dense-pack
  - Columns exhibit significantly less entropy
  - Examples:

    | Male, Female, Female, Female, Male |
    | 1998, 1998, 1999, 1999, 1999, 2000 |

  - Caveat: CPU cost (use lightweight compression)

Ack: Slide from  VLDB 2009 tutorial on Column store

# Telco example explained (3/3): sorting & indexing efficiency

1 Compression and dense-packing free up space

- 1 Use multiple overlapping column collections
- 1 Sorted columns compress better
- 1 Range queries are faster
- 1 Use sparse clustered indexes

Ack: Slide from  VLDB 2009 tutorial on Column store