# CSE 444: Database Internals

Lecture 20
Distributed Query Processing
and Optimization

# References

- In our book: Sections 20.3 and 20.4

- Other book: Database management systems.
  Ramakrishnan and Gehrke.
  Third Ed. Chapter 22

# Where We Are

- We know all about a single-machine DBMS

- What if we had more than one machine?
  – Distributed query processing (today)
  – Distributed transactions (next lecture)
  – Replication (next week)
  – Parallel query processing (later next week)

# Outline
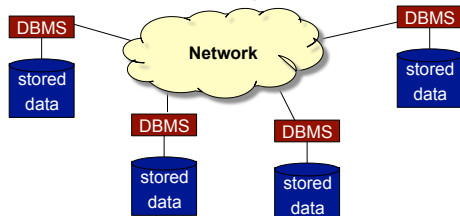
- Distributed DBMS motivation

- Distributed query optimization/processing

- Distributed DBMS limitations and challenges

# Distributed DBMS

- Important: many forms and definitions
- Our definition: shared nothing infrastructure
  – Multiple machines connected with a network
  – Machines can be far away from each other



# Vertical Partitioning

**Students**

| StudentID | Name | Address | GPA | Status |
|-----------|------|---------|-----|--------|
| 234234 | Mary | Houston | 3.5 | Scholarship |
| 345345 | Sue | Seattle | 2.9 | Fellowship |
| 345343 | Joan | Seattle | 3.2 | Veteran |
| 234234 | Ann | Portland | 3.9 | None |

**Registrar**

| StudentID | Name | Address | GPA |
|-----------|------|---------|-----|
| 234234 | Mary | Houston | 3.5 |
| 345345 | Sue | Seattle | 2.9 |
| … | … | … | … |

**Bursar**

| SID | Status |
|-----|--------|
| 234234 | Scholarship |
| 345345 | Fellowship |
|  |  |

1

## Horizontal Partitioning

**Customers**

| SSN | Name | City | Country |
|---|---|---|---|
| 234234 | Mary | Houston | USA |
| 345345 | Sue | Seattle | USA |
| 345343 | Joan | Seattle | USA |
| 234234 | Ann | Portland | USA |
| -- | Frank | Calgary | Canada |
| -- | Jean | Montreal | Canada |

**CustomersInHouston**

| SSN | Name | City | Country |
|---|---|---|---|
| 234234 | Mary | Houston | USA |

**CustomersInSeattle**

| SSN | Name | City | Country |
|---|---|---|---|
| 345345 | Sue | Seattle | USA |
| 345343 | Joan | Seattle | USA |

**CustomersInCanada**

| SSN | Name | City | Country |
|---|---|---|---|
| -- | Frank | Calgary | Canada |
| -- | Jean | Montreal | Canada |

Magda Balazinska - CSE 444, Spring 2012          7

---

## Reasons for a Distributed DBMS

- Scalability (ex: Amazon, eBay, Google)
  - Many small servers cheaper than one big server
  - Need to scale incrementally
  - Typically: horizontal partitioning
- Inherent distribution
  - Large organizations have data at multiple locations (different offices) -> original motivation
  - Different types of data in different DBMSs
  - Web-based and Internet-based applications
  - Either horizontal or vertical partitioning

Magda Balazinska - CSE 444, Spring 2012          8

---

## Distributed vs Parallel DBMS

- Parallel DBMS
  - Goal to improve performance through parallelization
    - Parallel load, parallel index build, parallel query processing
  - Data distribution is governed solely by performance
  - Typically all machines are on the same cluster
- Distributed DBMS
  - Data stored across several sites
  - Each site capable of running independent of other sites
  - Each site can be administered autonomously
  - Data distribution not only based on performance

Magda Balazinska - CSE 444, Spring 2012          9

---

## Goals of a Distributed DBMS

- Shield users from distribution details
- Distribution transparency
  - Naming transparency
  - Location transparency
  - Fragmentation transparency
  - Performance transparency
    - Distributed query optimizer ensures similar performance no matter where query is submitted
  - Schema change and transaction transparency
- Replication transparency
- and more... Magda Balazinska - CSE 444, Spring 2012          10

---

## Distributed DBMS Features

- 70's and 80's, three main prototypes:
  - SDD-1, distributed INGRES, and R*

- Main components of a distributed DBMS
  - Defining data placement and fragmentation
  - Distributed catalog
  - Distributed query optimization (today)
  - Distributed transactions (next lecture)
  - Managing replicas (next week)

Magda Balazinska - CSE 444, Spring 2012          11

---

## Outline

- Distributed DBMS motivation

- Distributed query optimization/processing

- Distributed DBMS limitations and challenges

Magda Balazinska - CSE 444, Spring 2012          12

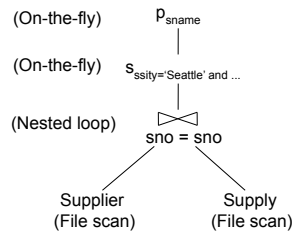## Review: Query Evaluation

**Steps of query evaluation**          **Query plan**

| Query Parser | (On-the-fly) | $p_{sname}$ |

(On-the-fly)  $s_{ssity='Seattle' \text{ and } ...}$

| Query Rewrite |

(Nested loop)  $\bowtie$
sno = sno

| Optimizer |

| Executor |          Supplier          Supply
                      (File scan)      (File scan)

---

## Review: Query Optimization

- Enumerate alternative plans
  - Many possible equivalent trees: e.g., join order
  - Many implementations for each operator

- Compute estimated cost of each plan
  - Compute number of I/Os and CPU utilization
  - Based on statistics

- Choose plan with lowest cost

---

## Distributed Query Optimization

- Search space is larger
  - Must select sites for joining relations
  - Must select method for shipping inner relation: whole or matches
- Minimize resource utilization
  - I/O, CPU, & *communication costs*
  - Example cost function used in R*
  - $W_{CPU} Nb_{inst} + W_{I/O} Nb_{I/O} + W_{msg} Nb_{msg} + W_{byte} Nb_{bytes}$
- Could also try to minimize response time
  - Least cost plan != Fastest plan

---

## Inner Table Transfer Strategy

- Ship whole
  - Read inner relation at its home site (using index or not)
  - Project inner relation to remove attributes not needed
  - Apply any single-table predicates
  - Ship results to site of outer relation and store in temp file
    - Note: we lose any indexes on the inner relation
- Fetch matches
  - For each tuple of outer, project tuple on join column
  - Send value to site of inner relation
  - Find matching tuples from inner relation
  - Ship projected, matching tuples back

---

## Additional Join Strategies

- Dynamically-created temporary index on inner
  - Ship inner relation, store in temp table, index temp
- Semijoin
  - Project outer relation on join column (eliminate dups)
  - Ship projected column to site with inner relation
  - Compute natural join and ship matching tuples back
  - Join the two relations
- Bloomjoin
  - Same idea as semijoin, but use Bloom filter instead of sending all values in the join column
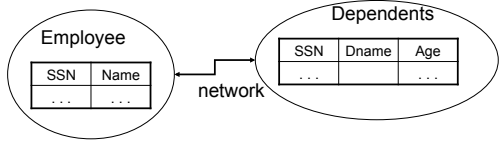  - Bloom filter creates some false positives through collisions

---

## Semijoin

- $R \ltimes S = \Pi_{A1,...,An} (R \bowtie S)$
- Where $A_1, \dots, A_n$ are the attributes in R
- Example:
  - Employee $\ltimes$ Dependents

## Semijoins in Distributed Databases



Employee

| SSN | Name |
|-----|------|
| ... | ... |

network

Dependents

| SSN | Dname | Age |
|-----|-------|-----|
| ... | | ... |

$Employee |\times|_{ssn=ssn} (s_{age>71} (Dependents))$

$T = \Pi_{SSN} (s_{age>71} (Dependents))$

$R = Employee \bowtie T$

$Answer = R \bowtie Dependents$

---

## Bloomjoin

Employee

| SSN | Name |
|-----|------|
| ... | ... |

Dependents

| SSN | Dname | Age |
|-----|-------|-----|
| ... | | ... |

bit vector: 1, 0, 1, 0, 0, 1, 1, 1

Ship bit vector

Hash(SSN)

Hash(SSN)

$R \bowtie Dependents$

Ship R containing only matching tuples

---

## Distributed vs Local Joins

Why can distributed joins be faster than local ones?

- More resources are available to the join
  - Ex: Distributed query can use twice the buffer pool (useful when accessing relations through unclustered indexes)

- Different parts of the join can proceed in parallel
  - Ex: Join tuples from page 1 while shipping page 2
  - Ex: Can sort the two relations in parallel

---

## Outline

- Distributed DBMS motivation

- Distributed query optimization/processing

- Distributed DBMS limitations and challenges

---

## Distributed DBMS Limitations

- Top-down
  - Global, a priori data placement
  - Global query optimization, one query at a time; no notion of load balance
  - Distributed transactions, tight coupling
- Assumes full cooperation of all sites
- Assumes uniform sites
- Assumes short-duration operations
- Limited scalability

---

## Distributed DBMS Challenges

In some distributed databases

- Autonomy: different administrative domains
  - Cannot always assume full cooperation
  - Do not want distributed transactions
- Heterogeneity
  - Different capabilities at different locations
  - Different data types, different semantics -> data integration pb
- Large-scale
  - Internet-scale query processor