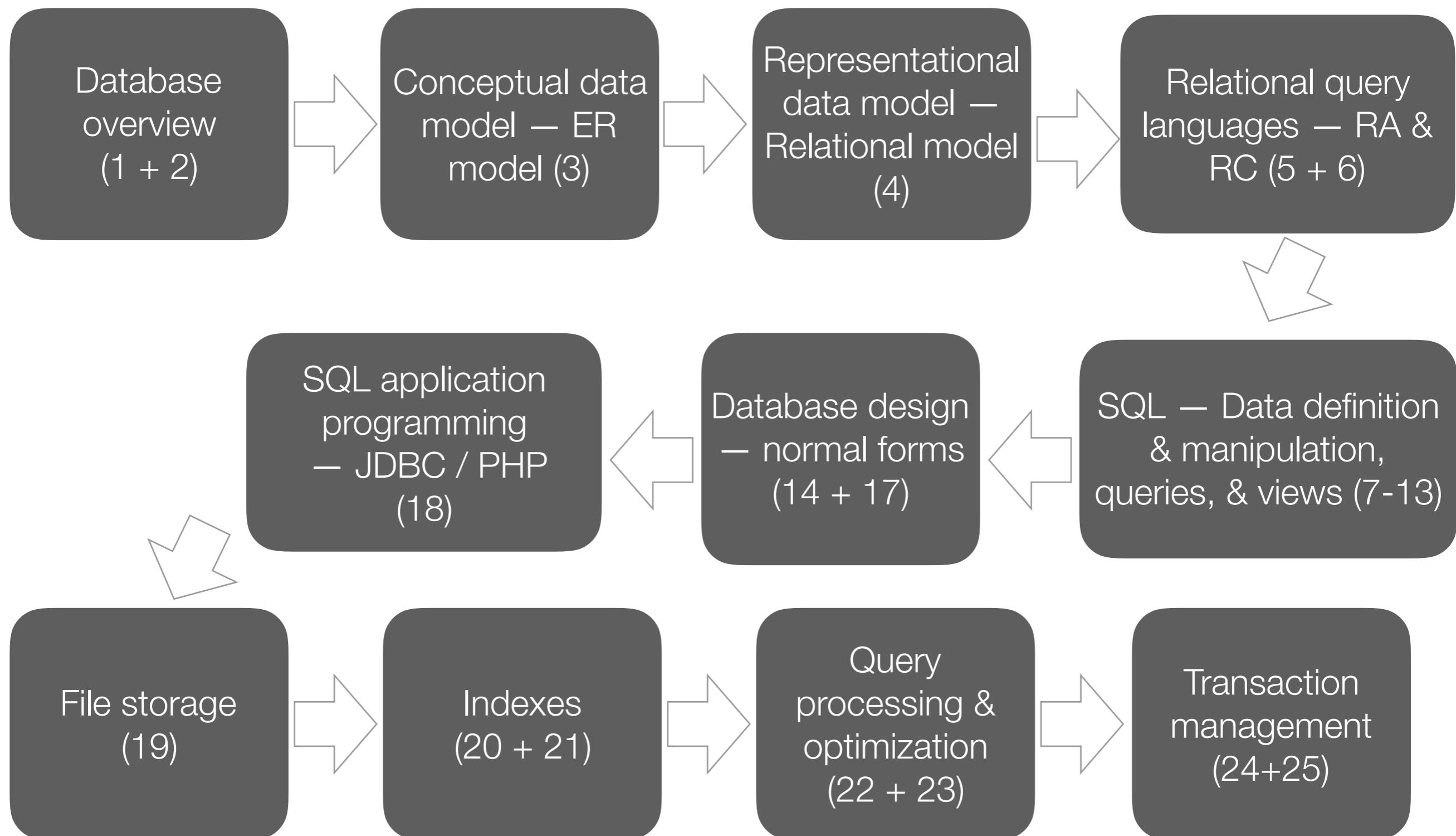


Big Data Systems

CS 377: Database Systems

Review: Course Material to Date



Review: What Has Been Covered

What I hope you've learned...

(1 + 2)

- Design a database

(Requirements -> ER diagram -> Relational model ->
Database normalization)

Database design
— normal forms

- Querying a database

(Relational algebra, calculus, SQL queries)

SQL application
programming
— JDBC / PHP
(12 & 22)

SQL – Data definition
& manipulation,
queries, & views (7-11)

File storage
(10)

Indexes
(17 + 18)

- Writing applications to use databases
(JDBC & SQL)

Query
processing &
optimization
(19 + 20)

Transaction
management
(21)

Review: What Has Been Covered

Database
(1 + 2)

Conceptual data
model (3)

Representational
data model –
Relational model
(4)

Relational query
languages – RA &
RC (5 + 6)

What I hope you've learned (at a high level)...

- Why some queries run faster on some systems compared to others
- How to think about optimizing your performance (Indexes, SQL Processing & optimization)

File structures
(16)

Indexes
(17 + 18)

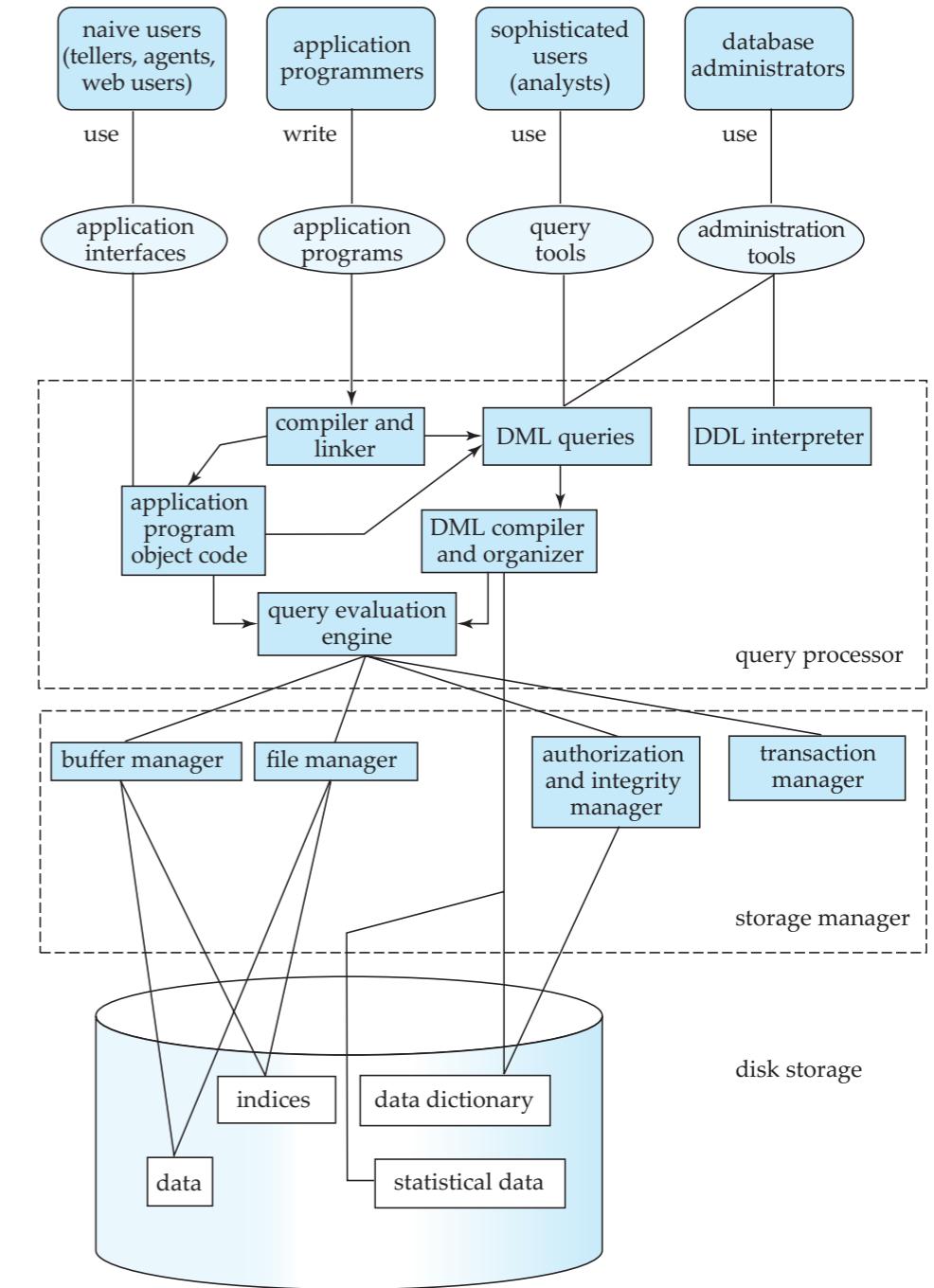
How to achieve ACID
(Logs & Locks)

Query
processing &
optimization
(19 + 20)

Transaction
management
(21)

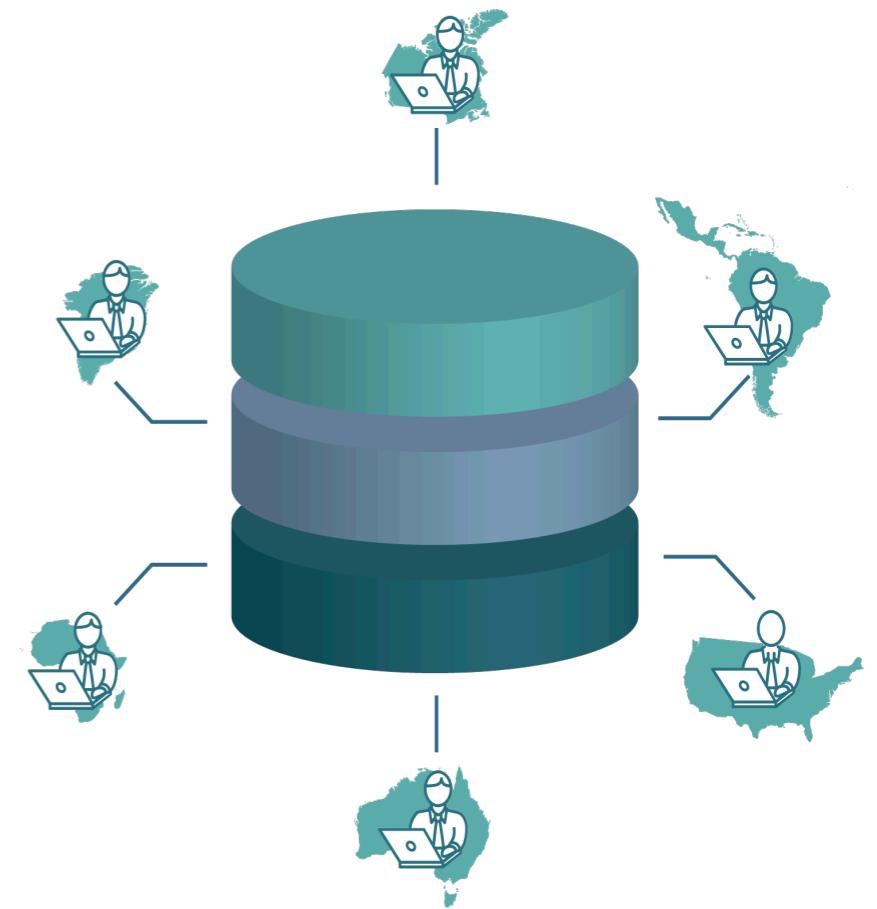
Review: “Peeking” Under the Hood

- Most aspects of traditional RDBMS is understood
- Learned enough to be “dangerous”
- Additional details can be picked up in courses or on your own

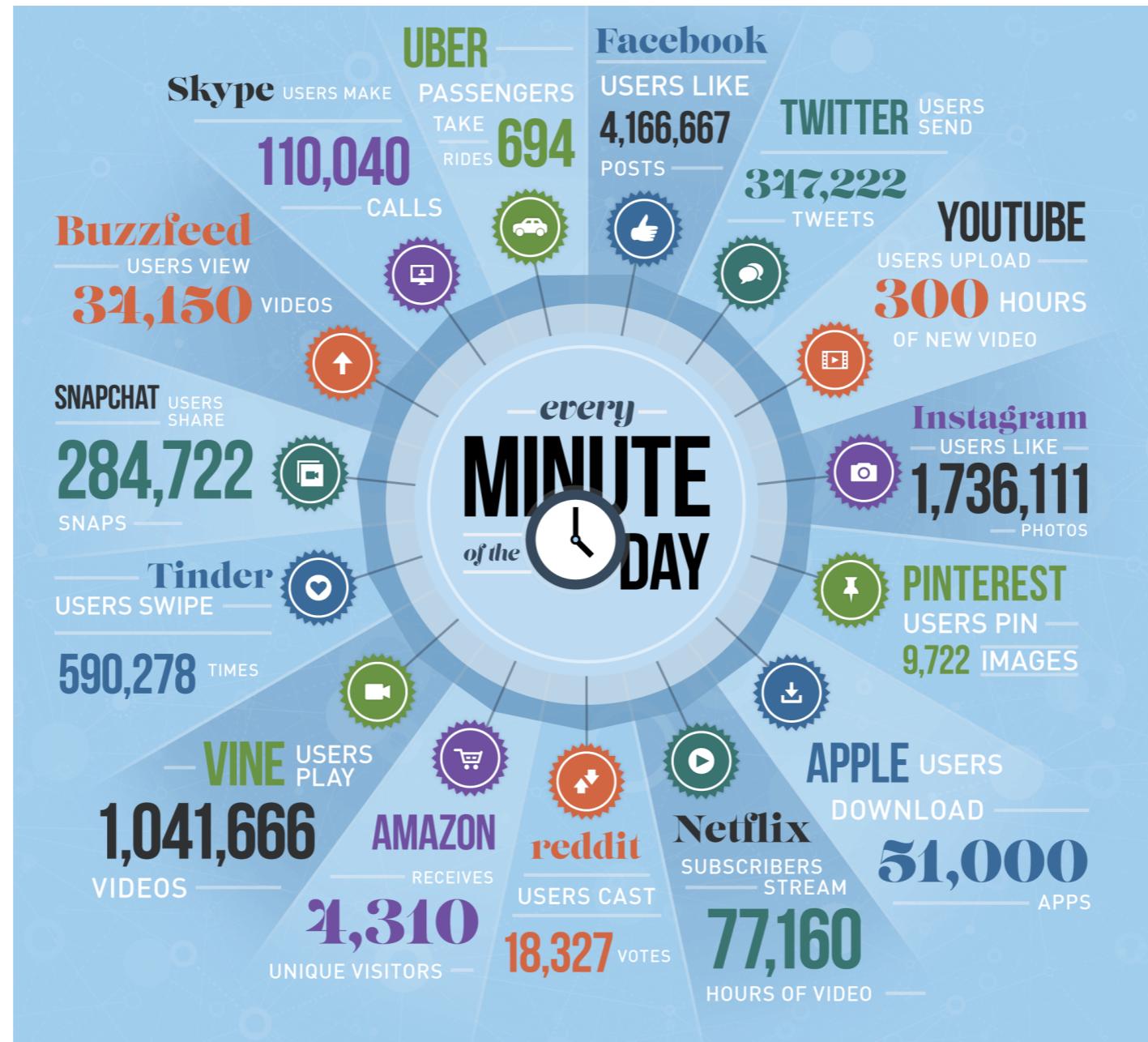


Review: Centralized Database

- Data is located in one place (one server)
 - All functions performed by the server
 - Query processing
 - Transaction management, concurrency control
 - ...
- What if I have 100 TB of data?



Data Never Sleeps

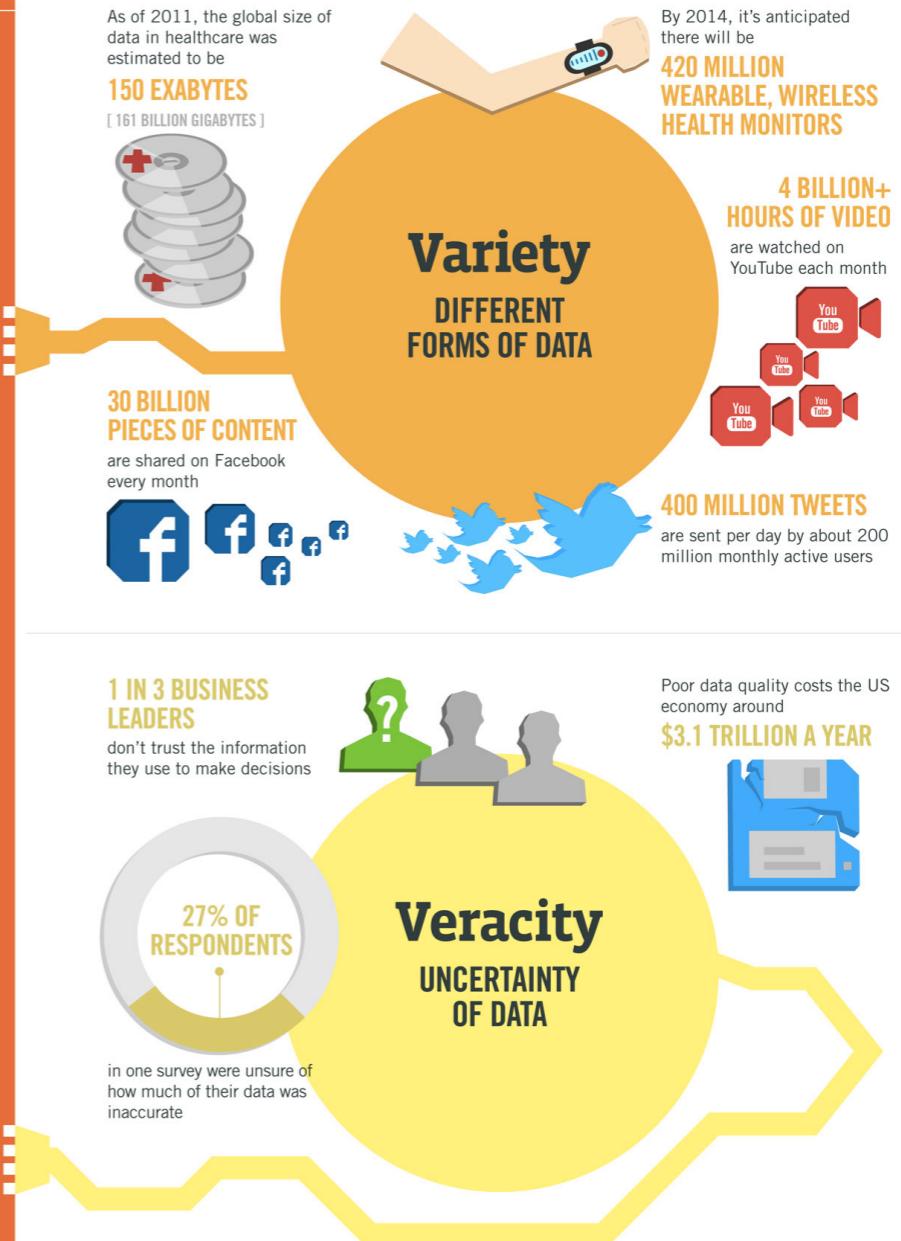
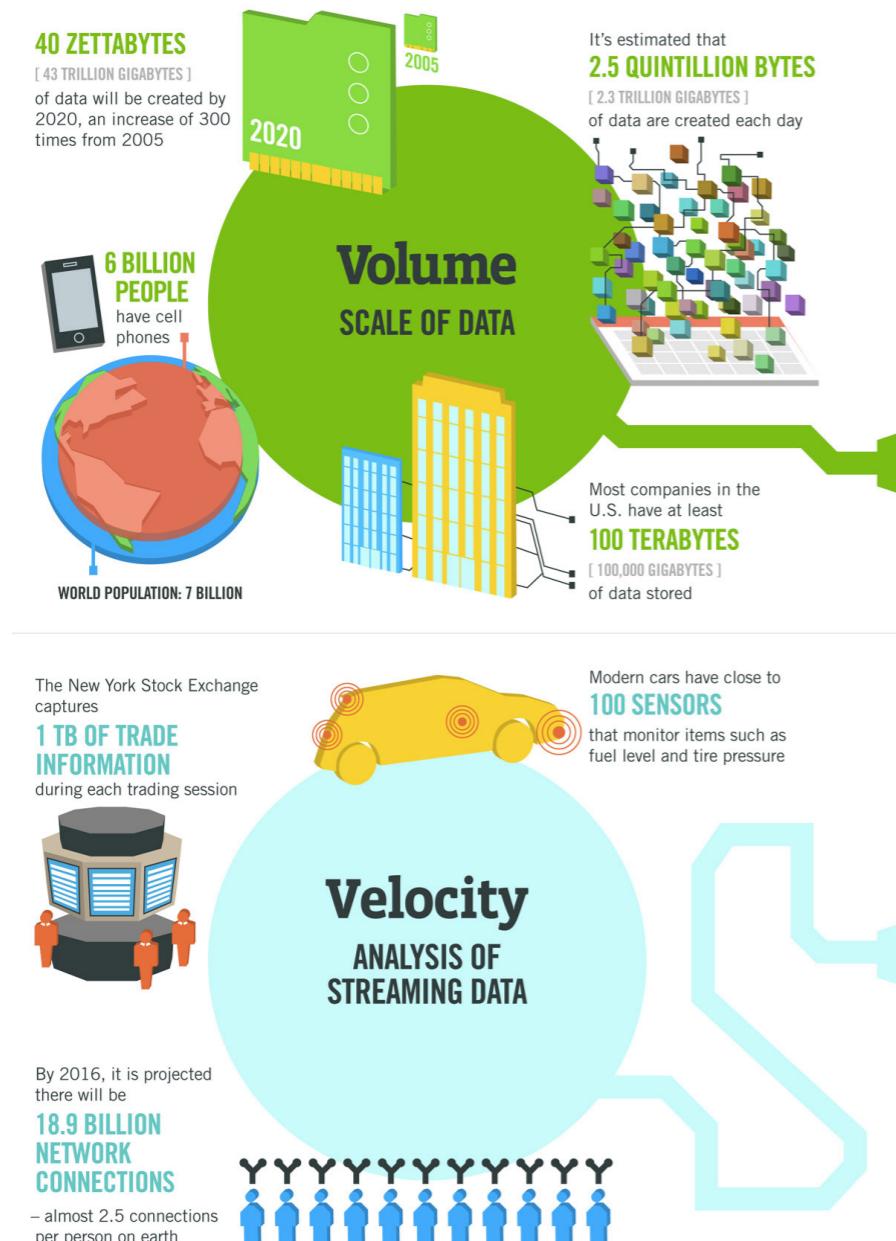


<https://www.domo.com/blog/2015/08/data-never-sleeps-3-0/>

Goal of Today's Lecture

- High-level overview of dealing with “big data”
 - What is big data?
 - What are different technologies I can use?
- Not meant to be detailed examination of all aspects of systems covered

4 V's of Big Data



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS



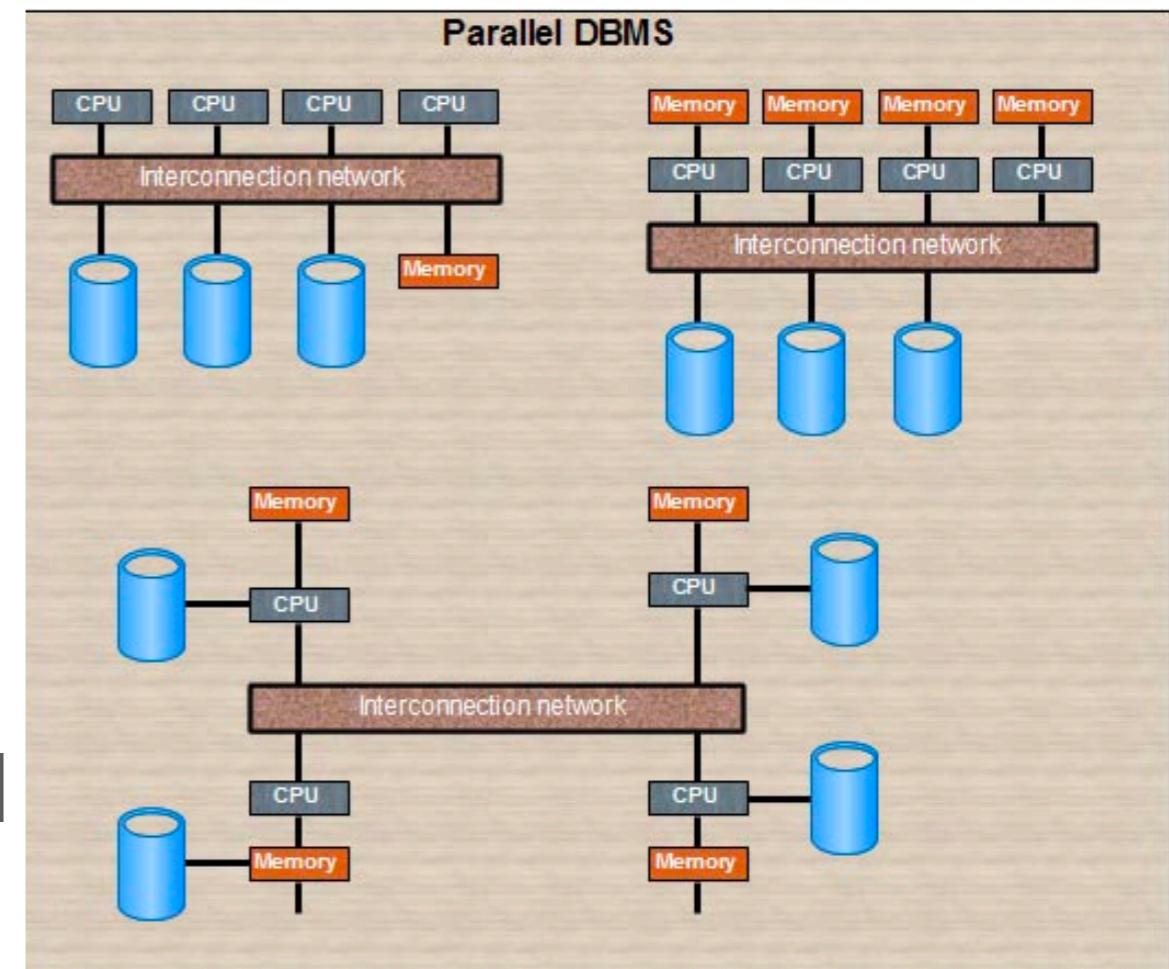
<http://www.ibmbigdatahub.com/infographic/four-vs-big-data>

Parallel & Distributed DBs: Motivation

- Single, monolithic DBMS is impractical and expensive
- Improve performance
- Increased availability & reliability
- Potentially lower cost of ownership
- Easier, more economical system expansion

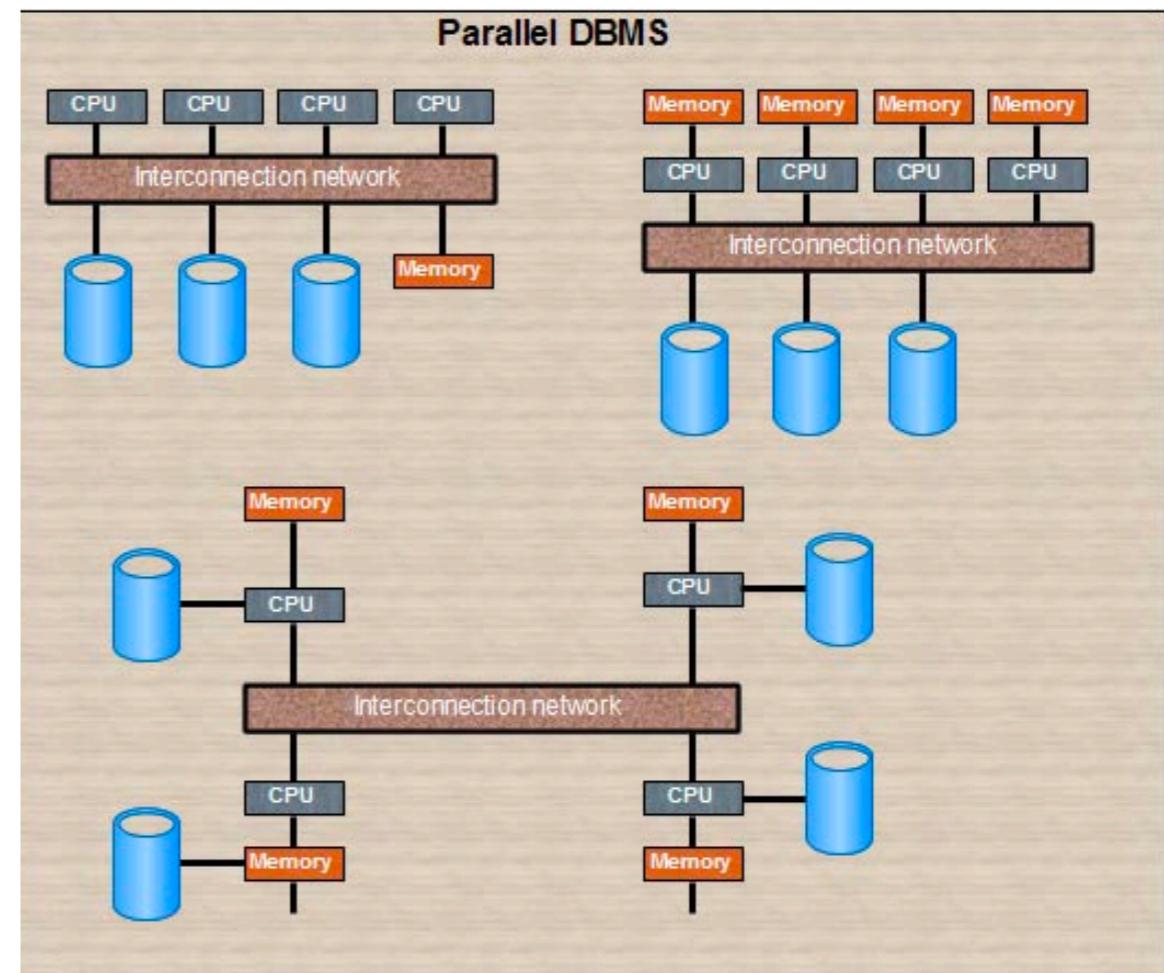
Parallel & Distributed DBs: Overview

- Data partitioned across multiple disks
- Allows parallel I/O for better speed-up
- Queries can be run in parallel with each other



Parallel & Distributed DBs: Overview

- Each processor can work independently on its own partition
- Individual relational operations (e.g., sort, join, aggregation) can be executed in parallel
- Concurrency control takes care of conflicts

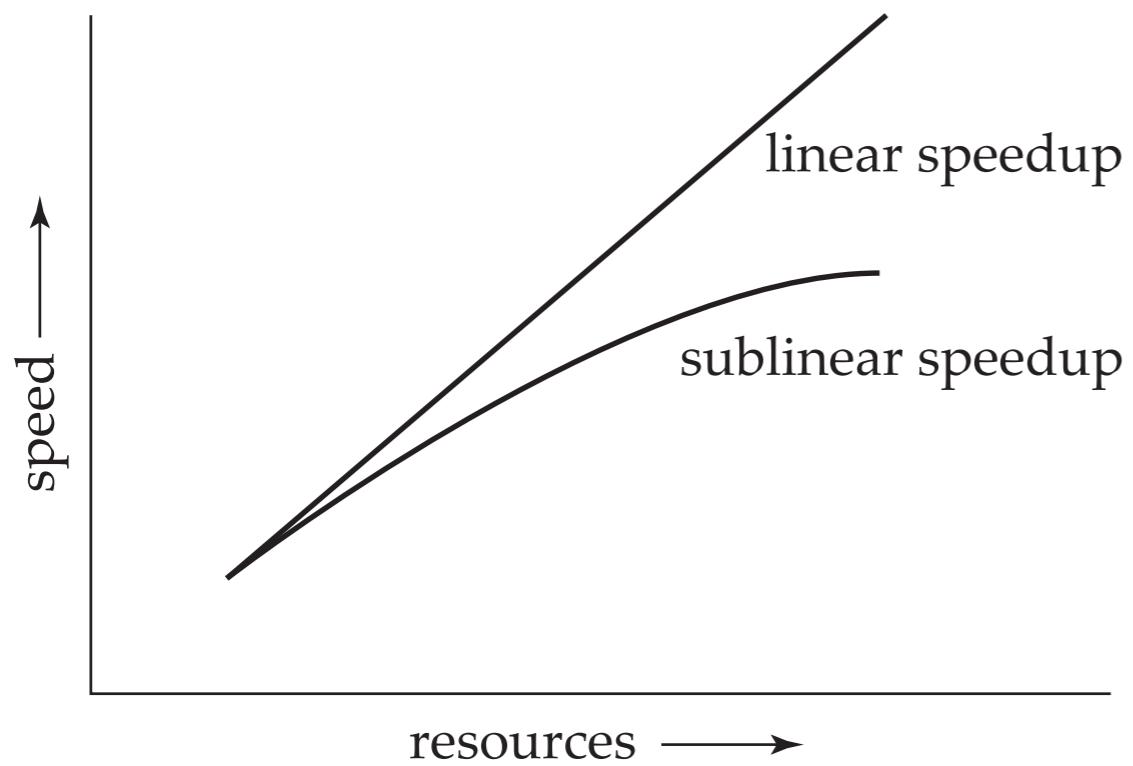


Scale-Up vs Scale-Out

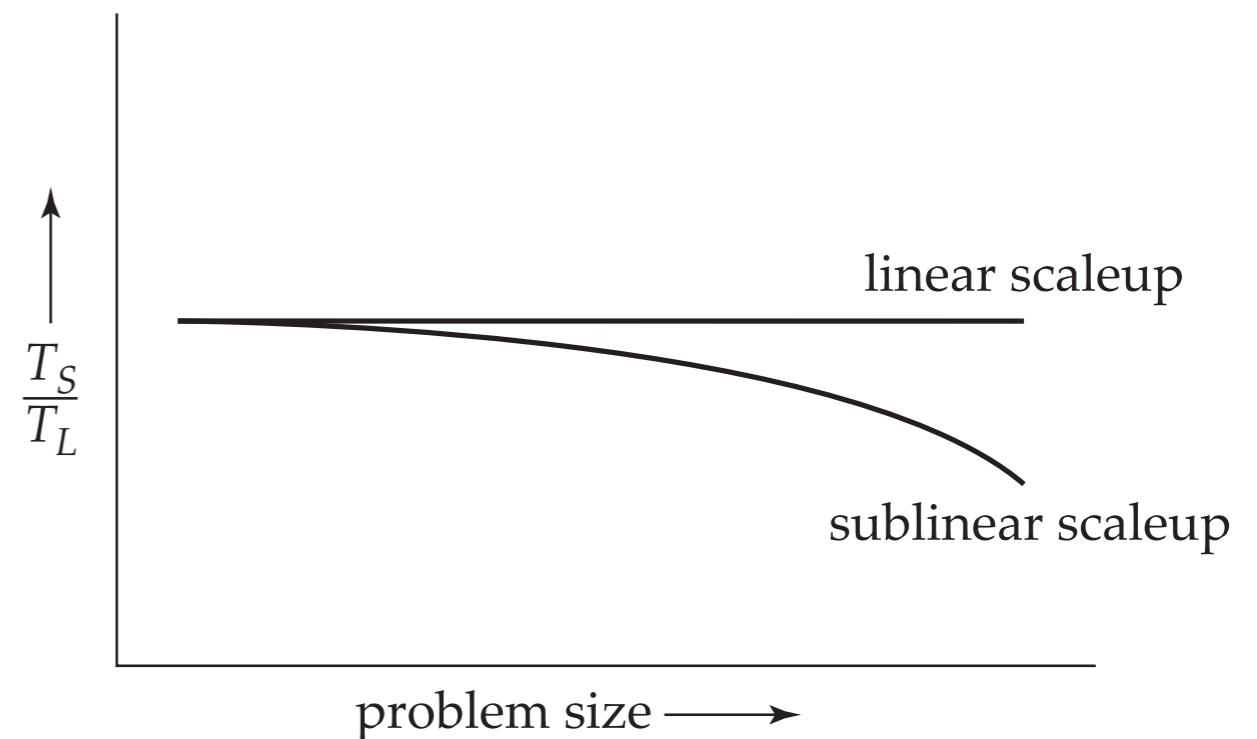
- Terminology to measure performance
- Speed-up: using more processors, how much faster will the task run (assuming same problem size)?
- Scale-up: using more processors, does performance remain the same as we increase problem size?
- Scale-out: using a larger number of servers, does performance improve?

Ideal Scenarios

Speed-up



Scale-up



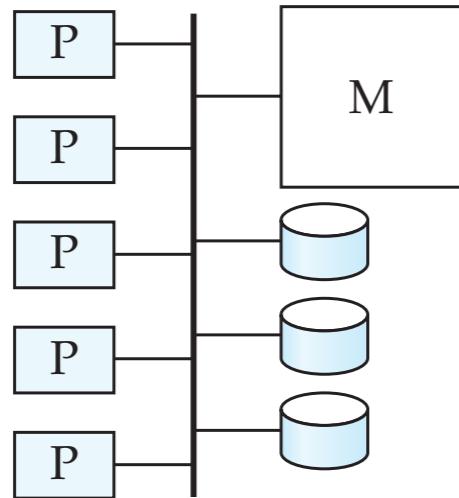
Parallel & Distributed DBs: Issues

- How to distribute the data
- How to optimize the cost of queries
 - Data transmission + local processing
- How to perform concurrency control
- How to make system resilient to failures and achieve atomicity & durability

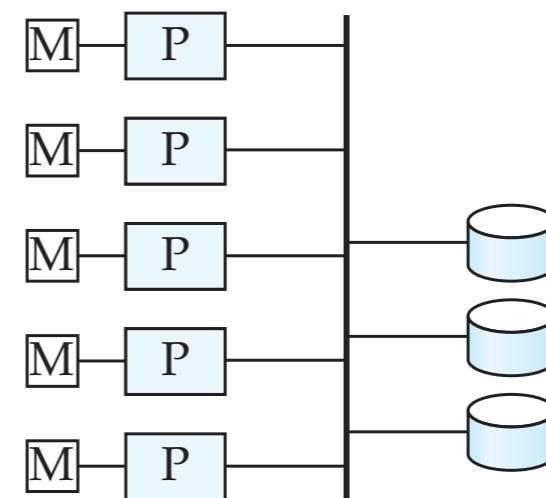
Parallel vs Distributed

- Parallel DBMS:
 - Nodes are physically close to each other
 - Nodes connected via high-speed LAN
 - Communication cost is small
- Distributed DBMS
 - Nodes can be far away
 - Nodes connected via public network
 - Communication cost and problems shouldn't be ignored

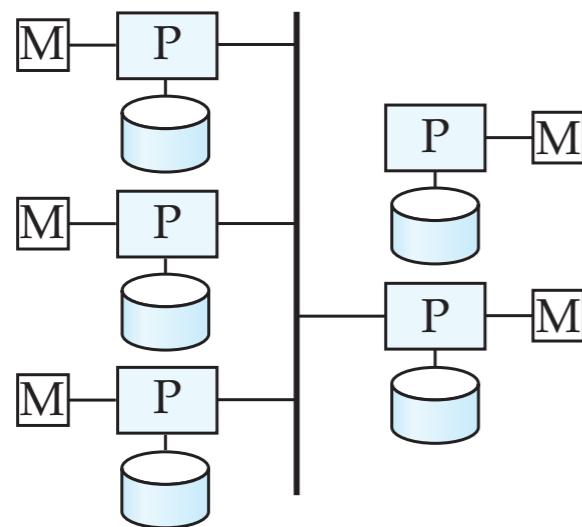
Parallel Architectures



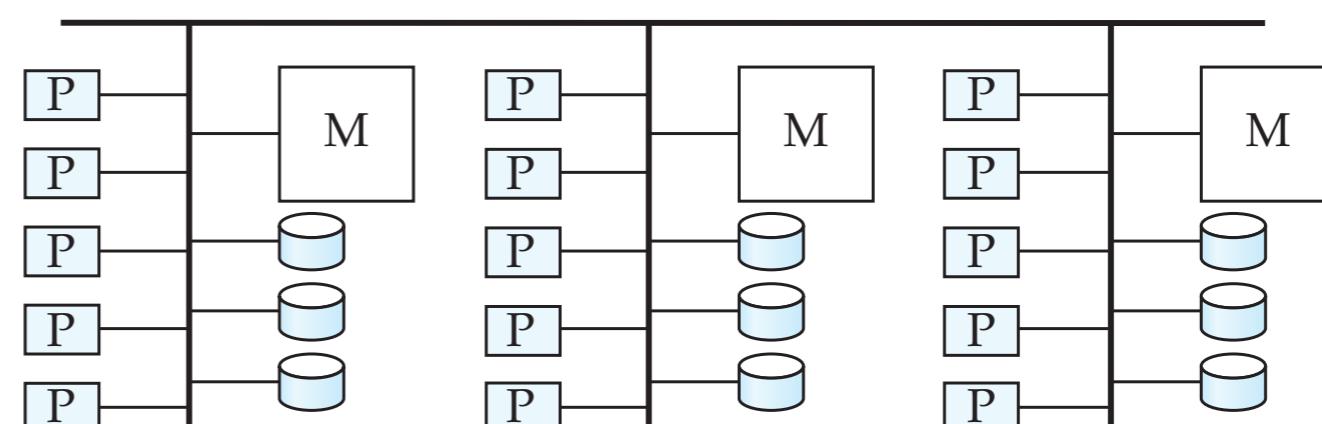
(a) shared memory



(b) shared disk



(c) shared nothing

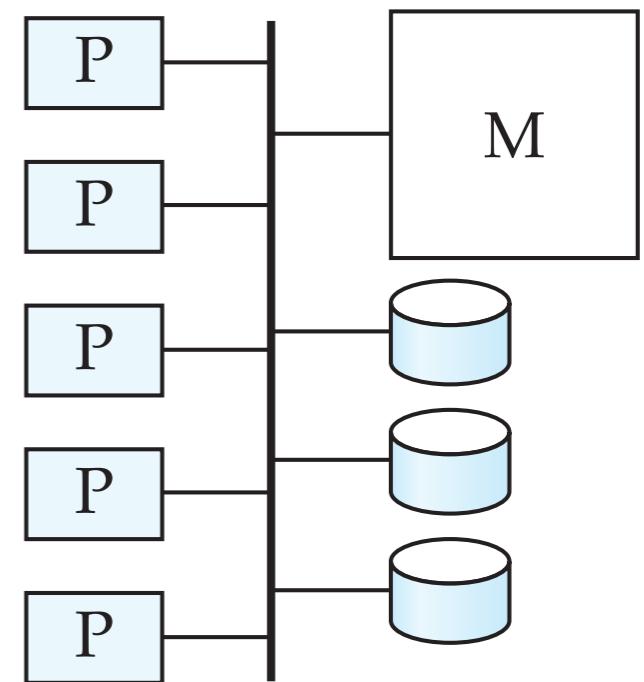


(d) hierarchical

Figure 17.8 (Database System Concepts)

Shared Memory

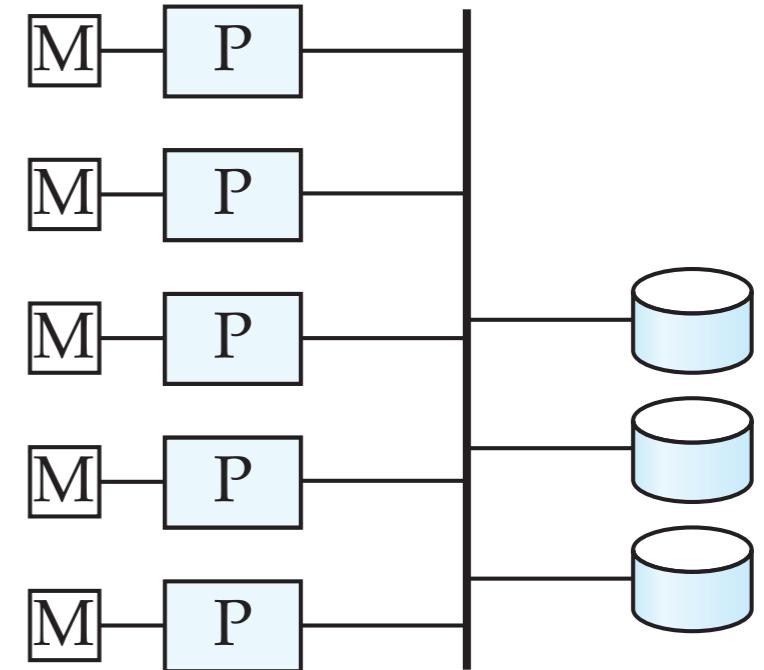
- Nodes share RAM + disk
- 10-100+ processors
- Easy to use and program
- Expensive to scale — last remaining cash cow in the hardware industry



(a) shared memory

Shared Disk

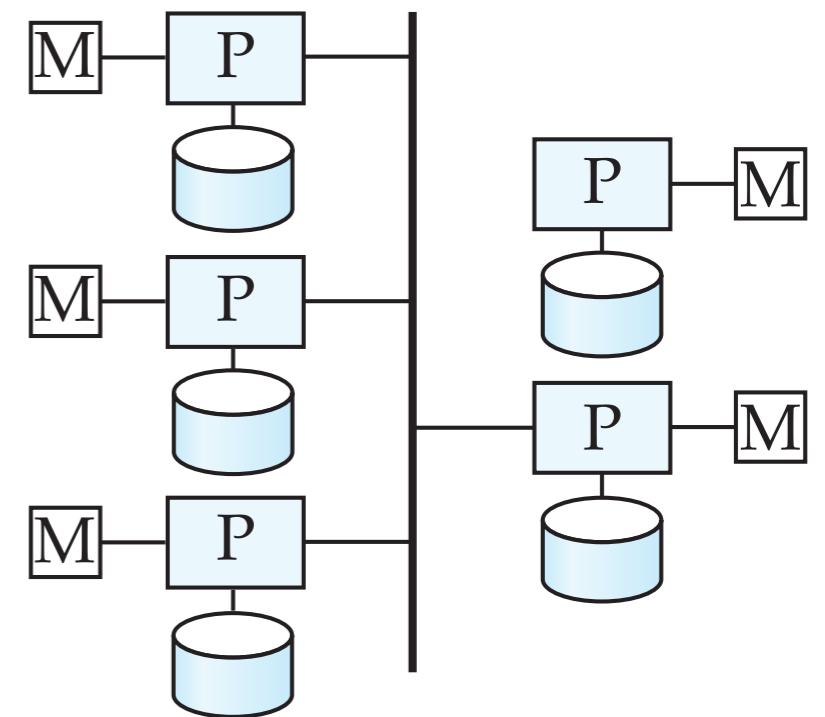
- Nodes share same disk
- Easy fault tolerance & consistency
- Hard to scale past a certain point
 - existing deployments typically have fewer than 10 machines
- Example: Oracle servers use this paradigm quite a bit



(b) shared disk

Shared Nothing

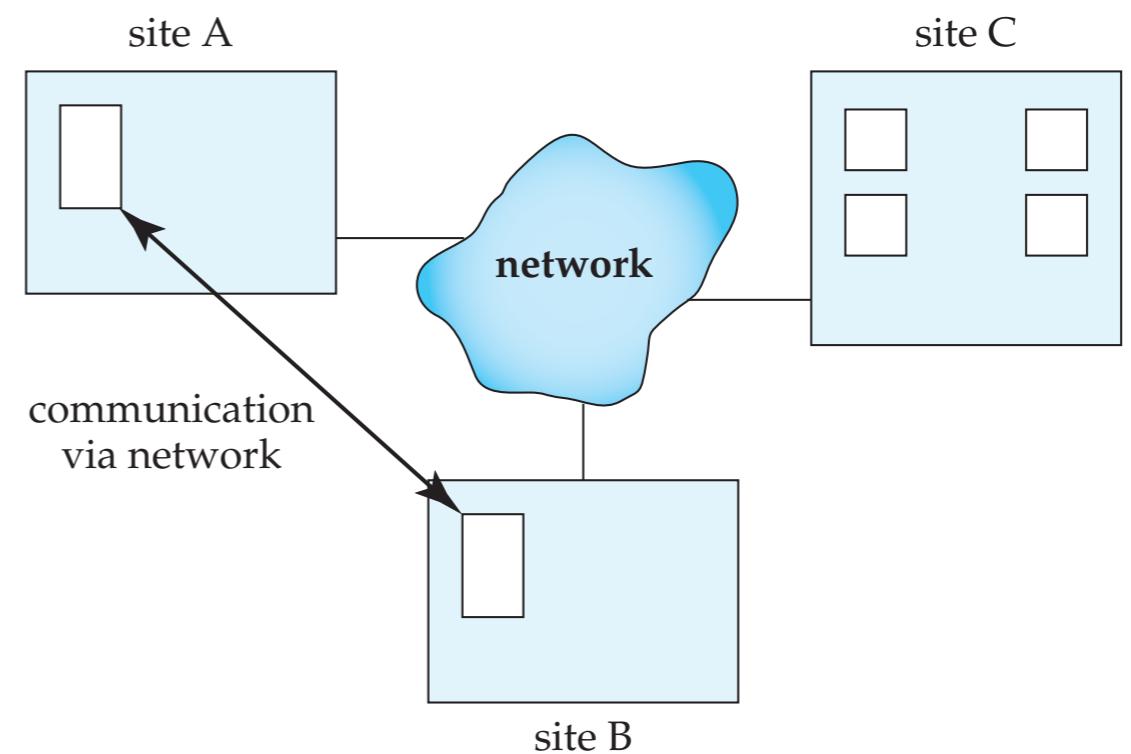
- Each instance has its own CPU, memory, and disk
- Easy to increase capacity
- Hard to ensure consistency
- Most scalable architecture but difficult to administer & tune



(c) shared nothing

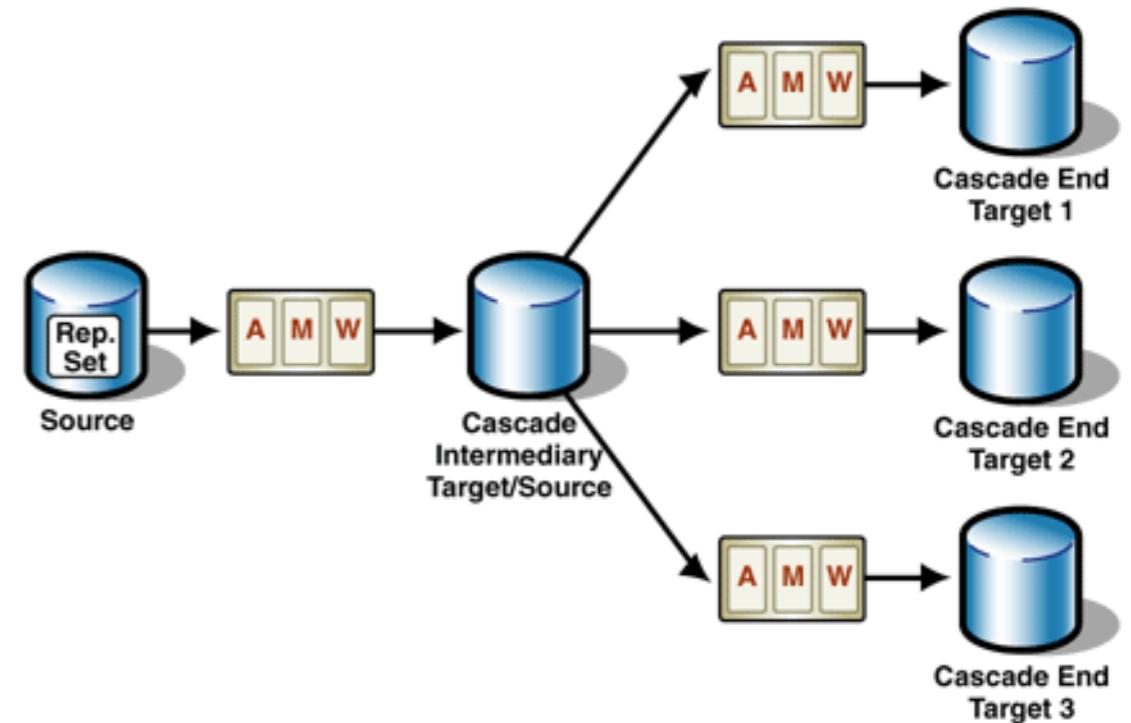
Distributed Databases

- Data spread over multiple machine
- Network interconnects the machines
- Similar to shared nothing architecture but larger communication cost



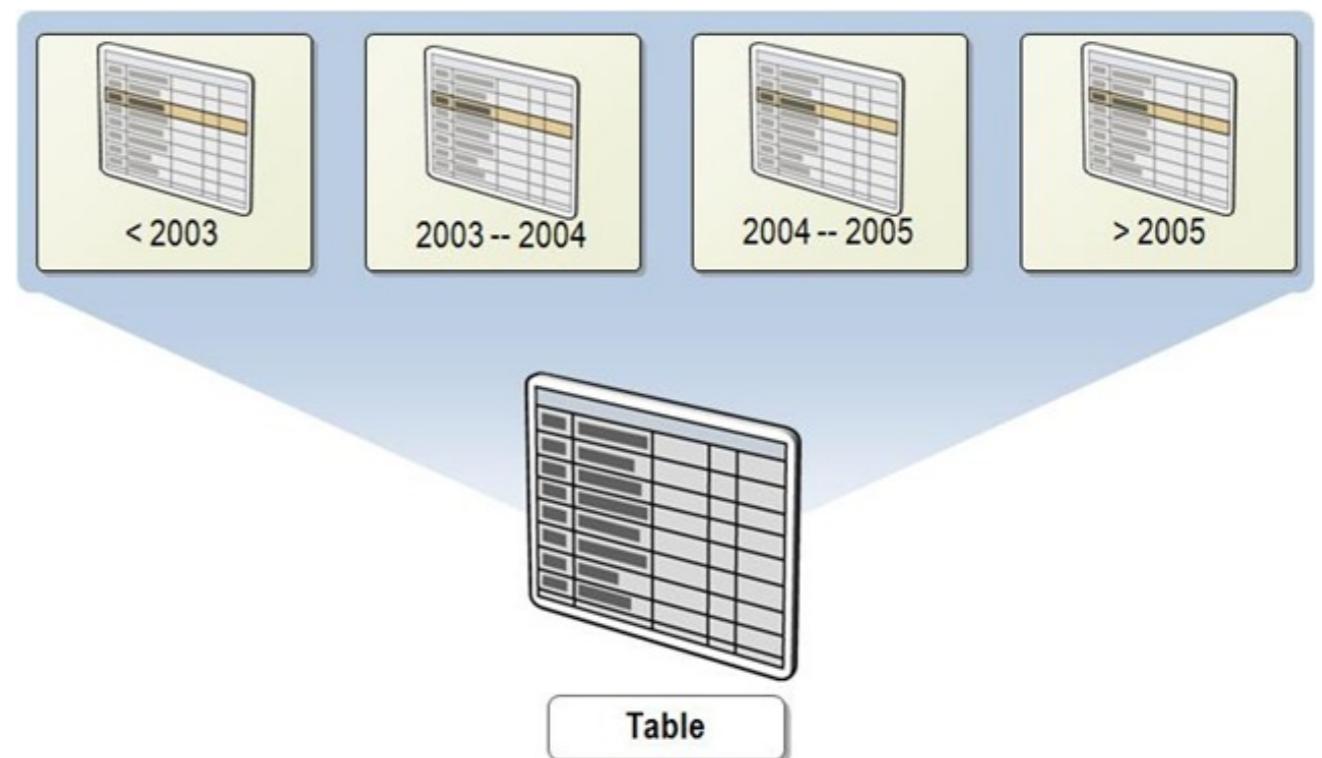
How to Distribute the Data?

- Replication: system maintains multiple copies of data
 - (PRO) Improves availability, parallelism, and reduced data transfer
 - (CON) Increased cost of updates, complexity of concurrency control



How to Distribute the Data?

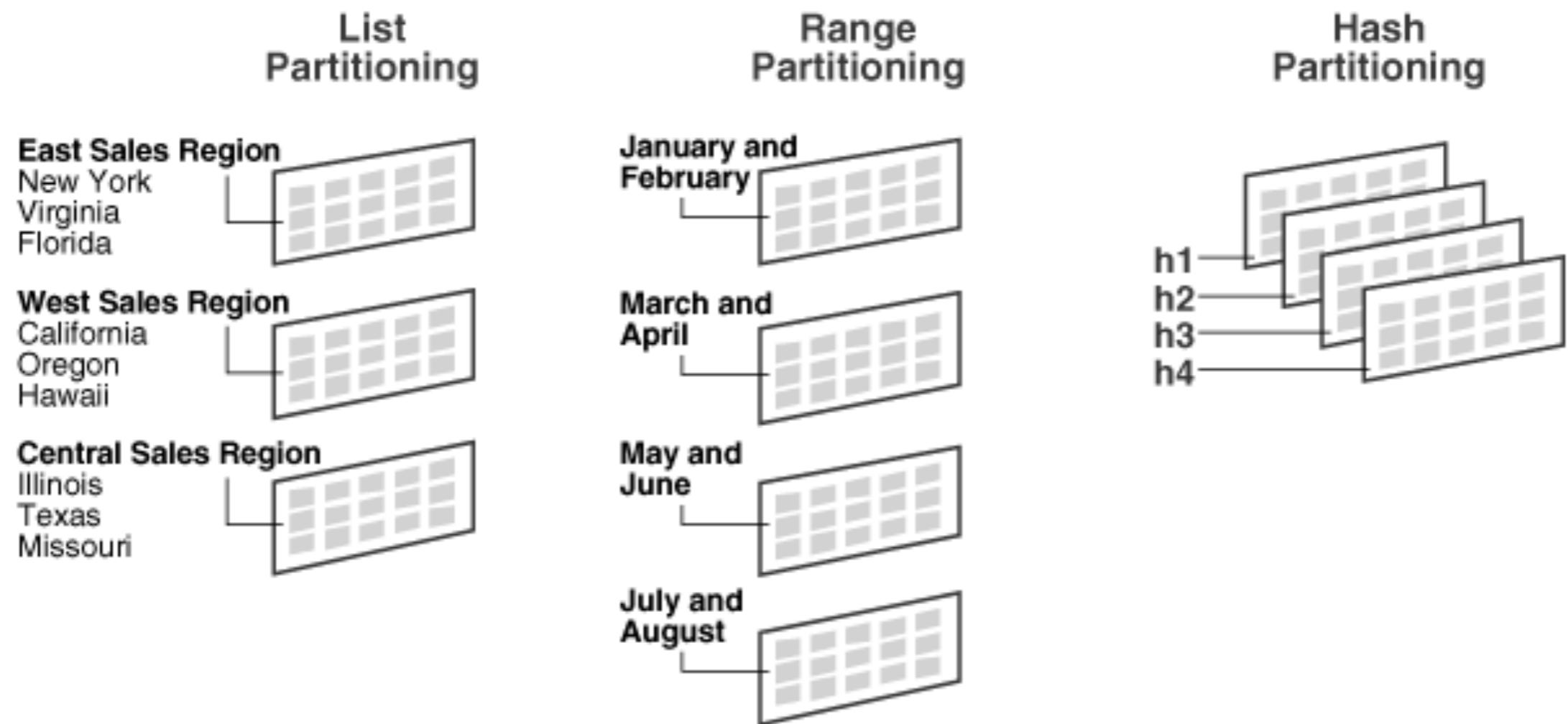
- Fragmentation: relation is partitioned into several fragments stored at distinct sites
- Combination of both replication & fragmentation



Fragmentation Strategies

- Horizontal partition: each tuple is assigned to one or more fragments
- Vertical partition: relation is split into smaller schemas each with a common candidate key to ensure lossless join

Horizontal Partition



https://docs.oracle.com/cd/B28359_01/server.111/b32024/partition.htm

Example: Horizontal Partition

<i>branch_name</i>	<i>account_number</i>	<i>balance</i>
Hillside	A-305	500
Hillside	A-226	336
Hillside	A-155	62

$account_1 = \sigma_{branch_name="Hillside"}(account)$

<i>branch_name</i>	<i>account_number</i>	<i>balance</i>
Valleyview	A-177	205
Valleyview	A-402	10000
Valleyview	A-408	1123
Valleyview	A-639	750

$account_2 = \sigma_{branch_name="Valleyview"}(account)$

Example: Vertical Partition

<i>branch_name</i>	<i>customer_name</i>	<i>tuple_id</i>
Hillside	Lowman	1
Hillside	Camp	2
Valleyview	Camp	3
Valleyview	Kahn	4
Hillside	Kahn	5
Valleyview	Kahn	6
Valleyview	Green	7

$deposit_1 = \Pi_{branch_name, customer_name, tuple_id} (employee_info)$

<i>account_number</i>	<i>balance</i>	<i>tuple_id</i>
A-305	500	1
A-226	336	2
A-177	205	3
A-402	10000	4
A-155	62	5
A-408	1123	6
A-639	750	7

$deposit_2 = \Pi_{account_number, balance, tuple_id} (employee_info)$

Example: Replication & Fragmentation

Figure 25.1

Data distribution and replication among distributed databases.

EMPLOYEES San Francisco
and Los Angeles

PROJECTS San Francisco

WORKS_ON San Francisco
employees

EMPLOYEES All

PROJECTS All

WORKS_ON All

EMPLOYEES New York

PROJECTS All

WORKS_ON New York
employees

EMPLOYEES Los Angeles

PROJECTS Los Angeles and
San Francisco

WORKS_ON Los Angeles
employees

EMPLOYEES Atlanta

PROJECTS Atlanta

WORKS_ON Atlanta
employees

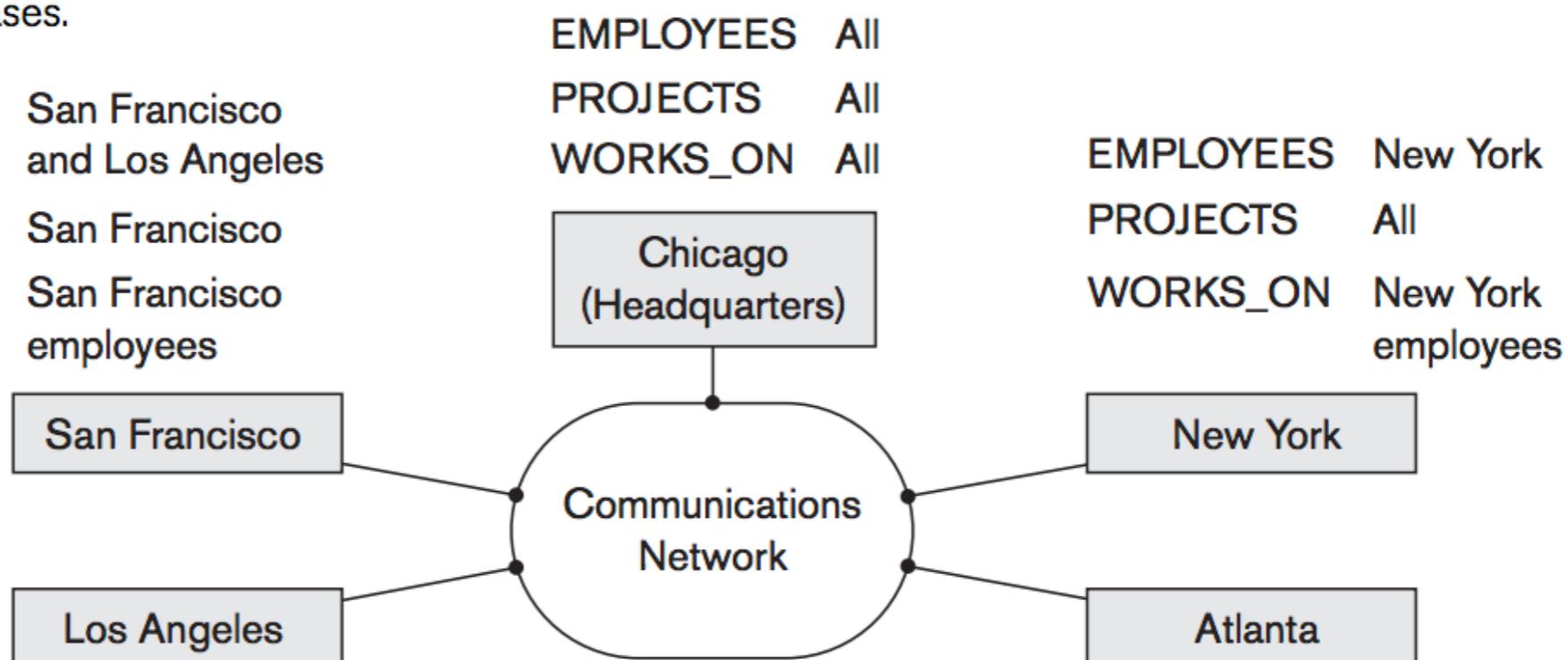


Figure 25.1 from FoDS book

Query Processing

- Single, centralized system – primary criterion for cost is just number of disk accesses
- Distributed system
 - Cost of data transmission over network
 - Potential gain in performance from having several sites process parts of the query

Review: Centralized DB Query

Given two relations $R(A, B)$ and $S(B,C)$ with no indexes, how do we compute the following?

- Selection: $\sigma_{A=123}(R)$

Linear search: scan file R and search for records A=123

Sort/hash for aggregation and apply sum

Review: Centralized DB Query

Given two relations $R(A, B)$ and $S(B, C)$ with no indexes, how do we compute the following?

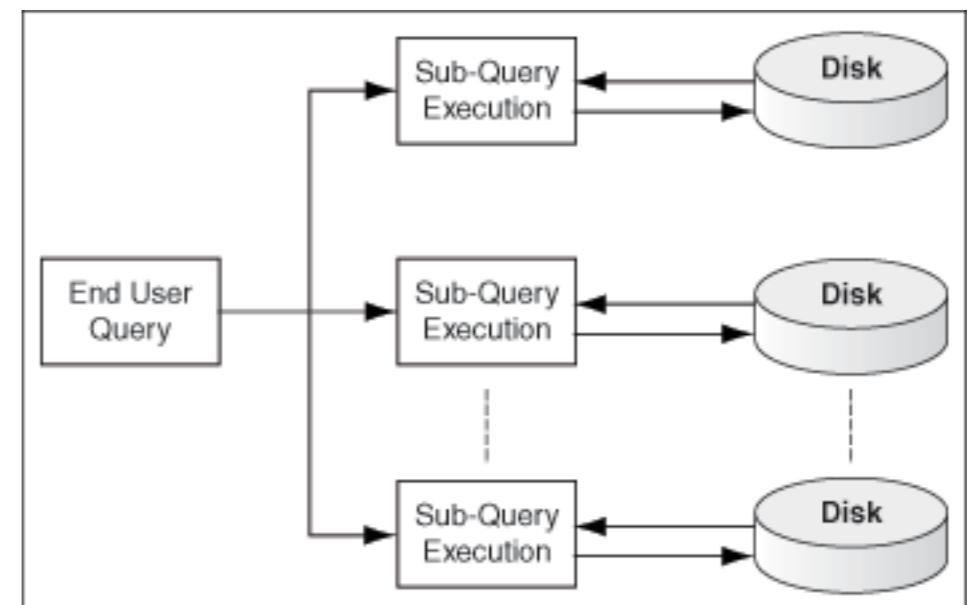
- Join: $R * S$
 - Nested block join
 - Hash join by creating hash index on B for smaller relation
 - Sort-merge join: sort on B for both relations

Parallel & Distributed DBs: Query

Given two relations $R(A, B)$ and $S(B, C)$ with horizontal partitioning and no indexes, how do we compute the following?

- Selection: $\sigma_{A=123}(R)$

Relatively straightforward – each machine scans its own partition and applies the condition

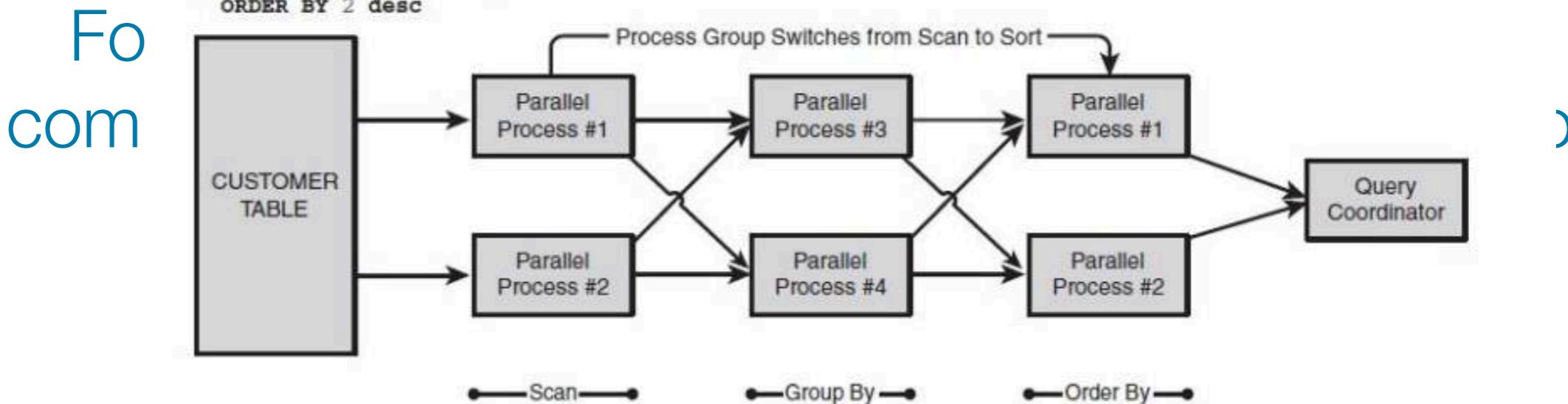


Parallel & Distributed DBs: Query

Given two relations $R(A, B)$ and $S(B, C)$ with horizontal partitioning and no indexes, how do we compute the following?

- Selection: $_A\mathcal{F}_{\text{SUM}(B)}(R)$

```
SELECT /*+ parallel(c,2) */ cust_last_name, count(*)
  FROM customers c
 GROUP BY cust_last_name
 ORDER BY 2 desc
```



Parallel & Distributed DBs: Query

Given two relations $R(A, B)$ and $S(B, C)$ with horizontal partitioning and no indexes, how do we compute the following?

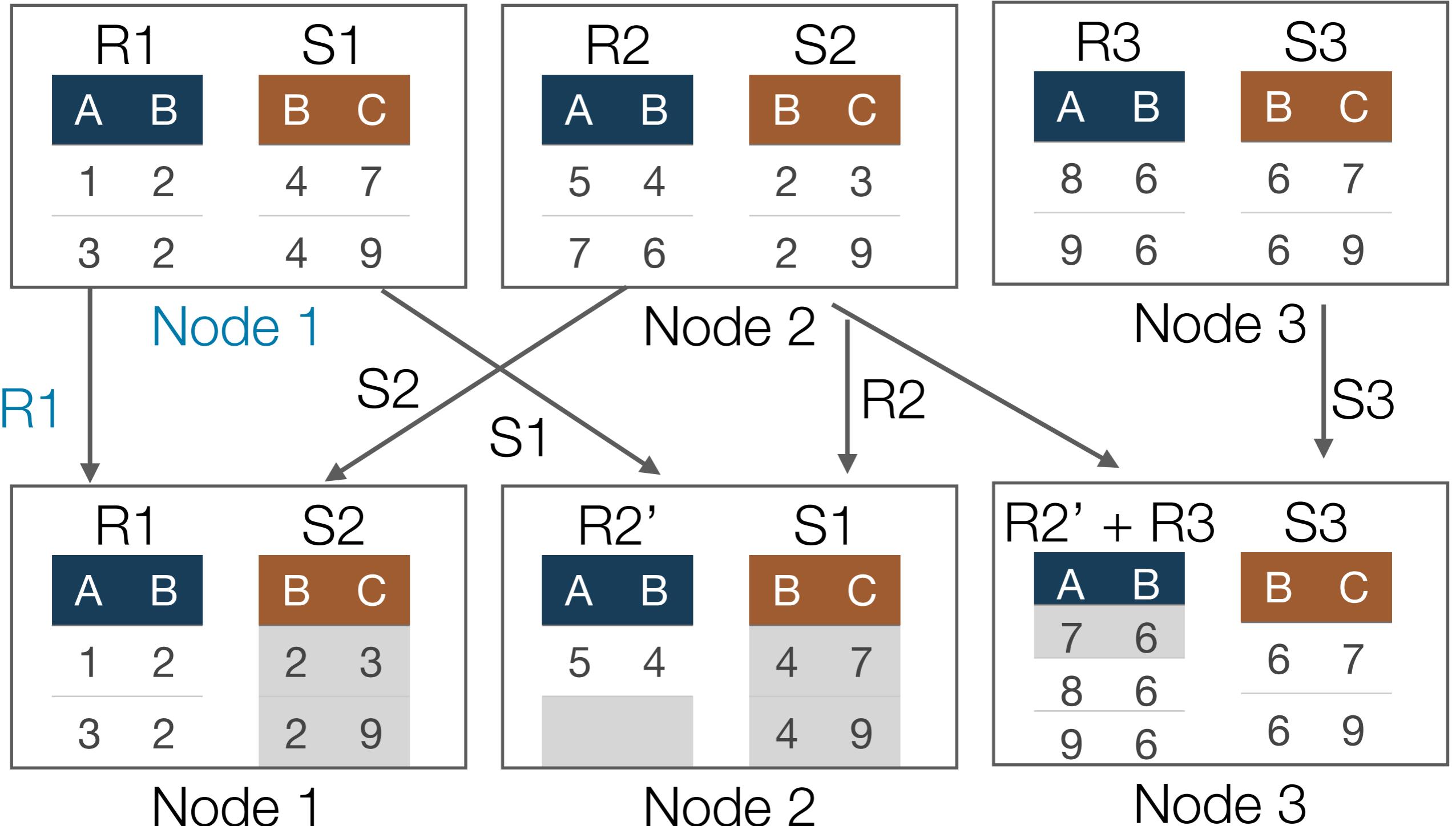
- Selection: ${}_A\mathcal{F}_{\text{SUM}(B)}(R)$

For hash and range partitions, relatively easy – complication occurs for round robin which needs to aggregate same values together

Parallel & Distributed DBs: Query

- Join: $R * S$
 - Strategy 1: Transfer both R and S into one central location and join (very expensive from sending)
 - Strategy 2: Perform local join by just sending the joining column of one relation, S, to where the other one is located, R (minimizes data transmission)

Example: Distributed Join



Example: Distributed Join (2)

A	B	C
1	2	3
1	2	9
3	2	3
3	2	9

Node 1

A	B	C
5	4	7
5	4	9

Node 2

A	B	C
7	6	7
7	6	9
8	6	7
...

Node 3

combine tuples for final output

Distributed Transactions & Recovery

- Dealing with multiple copies of data items — how to maintain consistency amongst the copies?
- Failure of individual sites — what to do when one site fails and then rejoins the system later?
- Failure of communication issues
- Distributed commit — what to do if some nodes fail during commit process?
- Distributed deadlock

Parallel & Distributed DBs: Properties

- Advantages
 - Data sharing
 - Reliability and availability
 - Improved query processing speed
- Disadvantages
 - May increase processing overhead
 - Harder to ensure ACID guarantees
 - More database design issues

What if I'm looking to manipulate diverse data?
For example, what if I want to extract links from
webpages and aggregate them by target
document? Should I do this all in SQL?

MapReduce



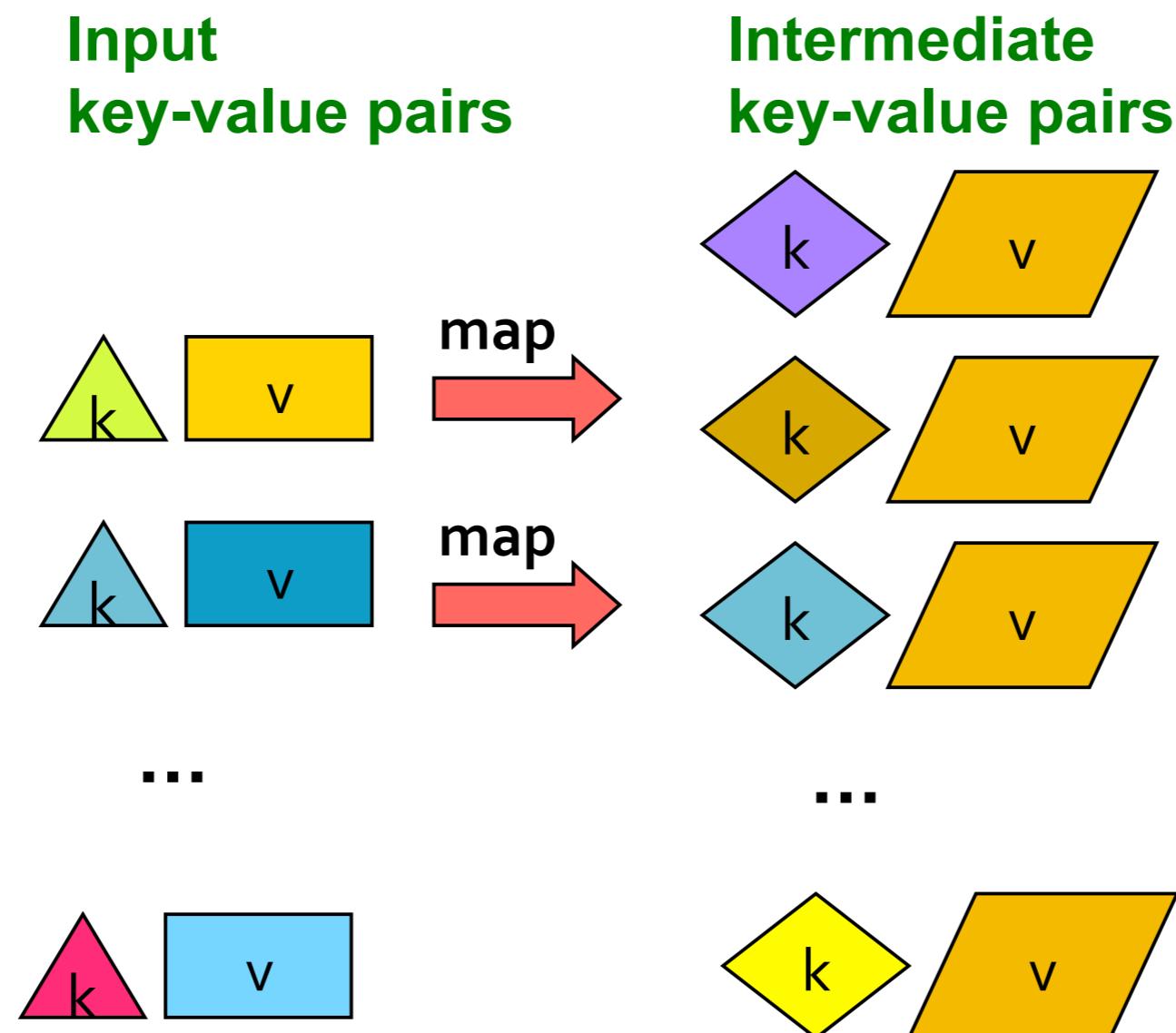
- Initially developed by Jeffrey Dean & Sanjay Ghemawat at Google [2004]
- Open source implementation: Apache Hadoop
- High-level programming model and implementation for large-scale parallel data processing
- Designed to simplify the task of writing parallel programs

MapReduce: Overview

- Read partitioned data
- Map: extract something you care about from each record
- Group by key: sort and shuffle (done by the system)
- Reduce: aggregate, summarize, filter, or transform
- Write the result

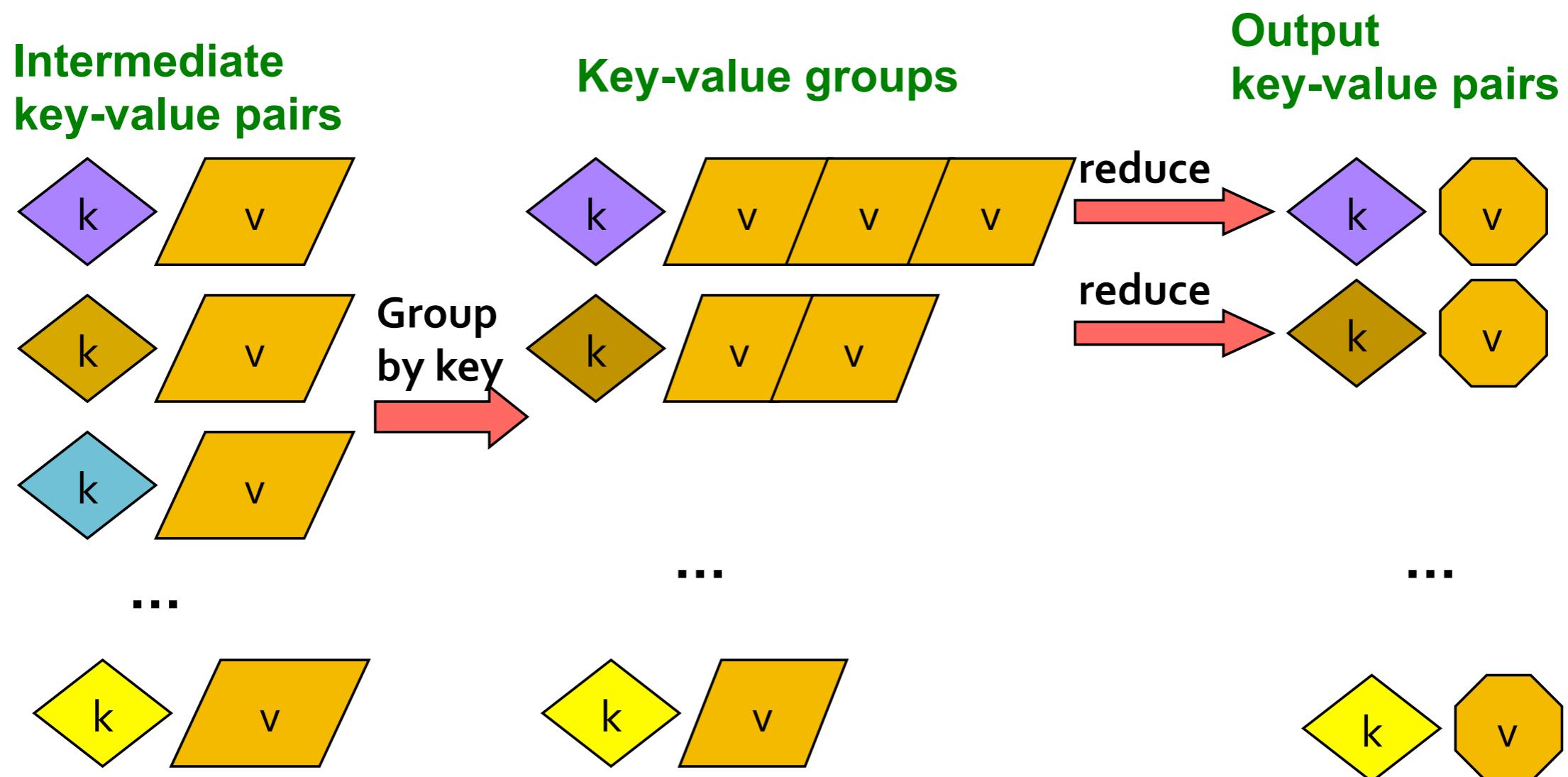
Outline stays the same, map and reduce should be tailored to the problem

MapReduce: Map Step



<http://www.mmds.org/#book>

MapReduce: Reduce Step



<http://www.mmds.org/#book>

Example: Word Counting

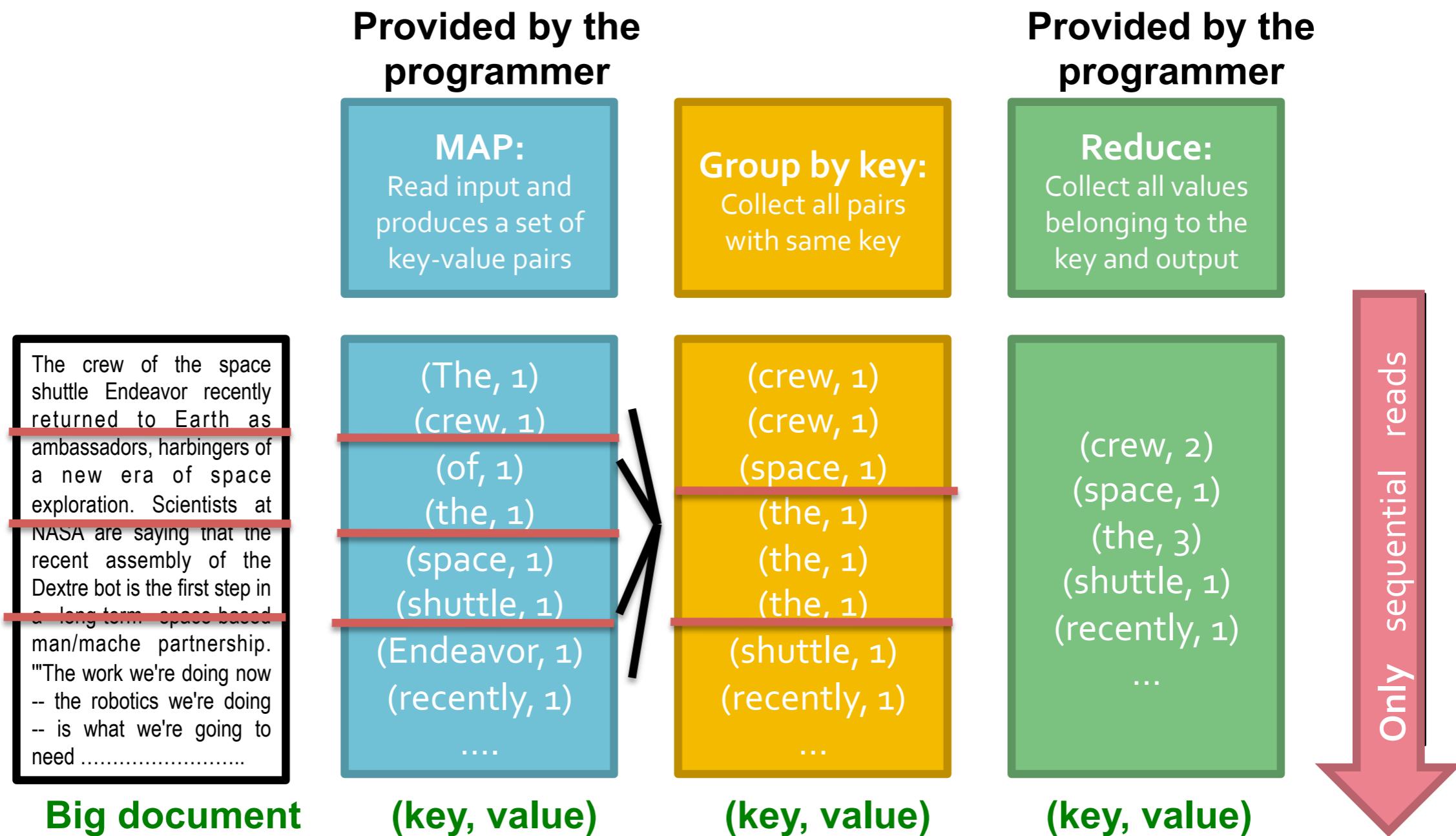
- We have a huge text document (~ 1 million words)
- Task: Count the number of times each distinct word appears in the file

Call me Ishmael. Some years ago--never mind how long precisely--having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world. It is a way I have of driving off the spleen and regulating the circulation. Whenever I find myself growing grim about the mouth; whenever it is a damp, drizzling November in my soul; whenever I find myself involuntarily pausing before coffin warehouses, and bringing up the rear of every funeral I meet; and especially whenever my hypos get such an upper hand of me, that it requires a strong moral principle to prevent me from deliberately stepping into the street, and methodically knocking people's hats off--then, I account it high time to get to sea as soon as I can. This is my substitute for pistol and ball. With a philosophical flourish Cato throws himself upon his sword; I quietly take to the ship. There is nothing surprising in this. If they but knew it, almost all men in their degree, some time or other, cherish very nearly the same feelings towards the ocean with me.

Example: Word Counting

- Traditional DBMS
 - Load document words into a table
 - SQL query:
**SELECT count(*)
FROM document
GROUP BY word**

Example: Word Counting (MapReduce)



<http://www.mmds.org/#book>

MapReduce Ecosystem

Many extensions to address limitations

- Capabilities to write directed acyclic graphs of MapReduce jobs (e.g., PIG by Yahoo!)
- Declarative languages (e.g., Hive by Facebook or SQL/Tenzing by Google)
- Increased integration of DBMS with MapReduce

Parallel DBMS vs MapReduce

Parallel DBMS

- Relational data model and schema
- Declarative query language (SQL)
- Easily combine operators into complex queries
- Query optimization, indexing, and physical tuning
- Streams data from one operator to next without blocking

Parallel DBMS vs MapReduce

MapReduce

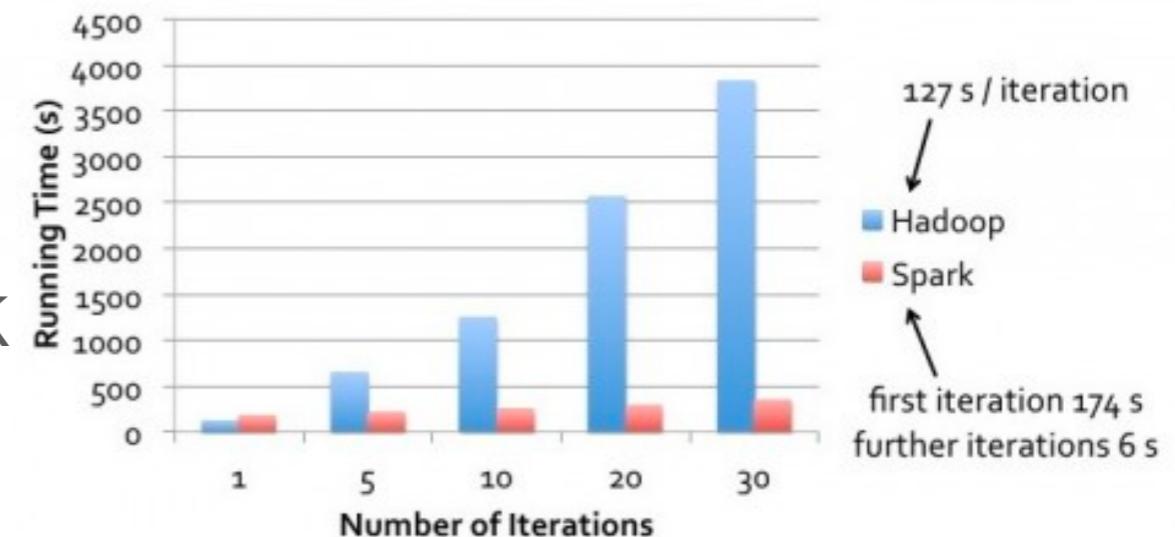
- Data model is file with key-value pairs
- Pre-loading data is not necessary before processing
- Easy to write user-defined operators
- Easily add nodes to the cluster
- Arguably more scalable, but also needs more nodes

Spark: MapReduce Replacement

- Tagline: Lightning-fast cluster computing
- Run programs up to 100x faster than MapReduce in memory or 10x faster on disk
- Easy to use with support for Java, Scala, Python, and R



Logistic Regression Performance



Big Data Systems: Recap

- Big Data (4 V's)
- Parallel/Distributed DBMS
 - Different architectures
 - Data distribution
 - Query processing
- MapReduce

