

Introduction to Data Structure (Data Management)

Lecture 15 – Parallel Databases (Ch. 20.1)

Felipe P. Vista IV



Chonbuk National University

- 1 -

Global Frontier College

Reminder

- Everybody, make sure that your name in ZOOM is in the following format:
 - University ID Num Name (no “()”)
 - Ex: 202054321 Juan Dela Cruz



- Not changing your name to this format
 - you might be marked Absent
 - * → absent?



- Why Parallel Processing
- Architectures
- Distributed Query Processing



INTRO TO DATA STRUCTURE

Why Compute in Parallel?

Why Compute in Parallel?

- **Multi-cores:**
 - Most processors have multiple cores
 - This trend will increase in the future
- **Big data:** too large to fit in main memory
 - Distributed query processing on 100-1000 servers
 - Widely available now using cloud services



Big Data

- Companies, org's, scientists have data that is **too big**
 - sometimes too complex,
 - to be managed without changing tools and processes
- Complex data processing:
 - Decision support queries (SQL w/ aggregates)
 - Machine learning (adds linear algebra and iteration)



Two Kinds of Parallel Data Processing

- Parallel databases, developed starting with the 80s
 - **OLTP** (Online Transaction Processing)
 - **OLAP** (Online Analytic Processing, or Decision Support)
- General purpose distributed processing:
MapReduce, Spark
 - **Mostly for Decision Support Queries**



Performance Metrics for Parallel DBMSs

P = the number of nodes (processors, computers)

- **Speedup:**

- More nodes, same data → higher speed

- **Scaleup:**

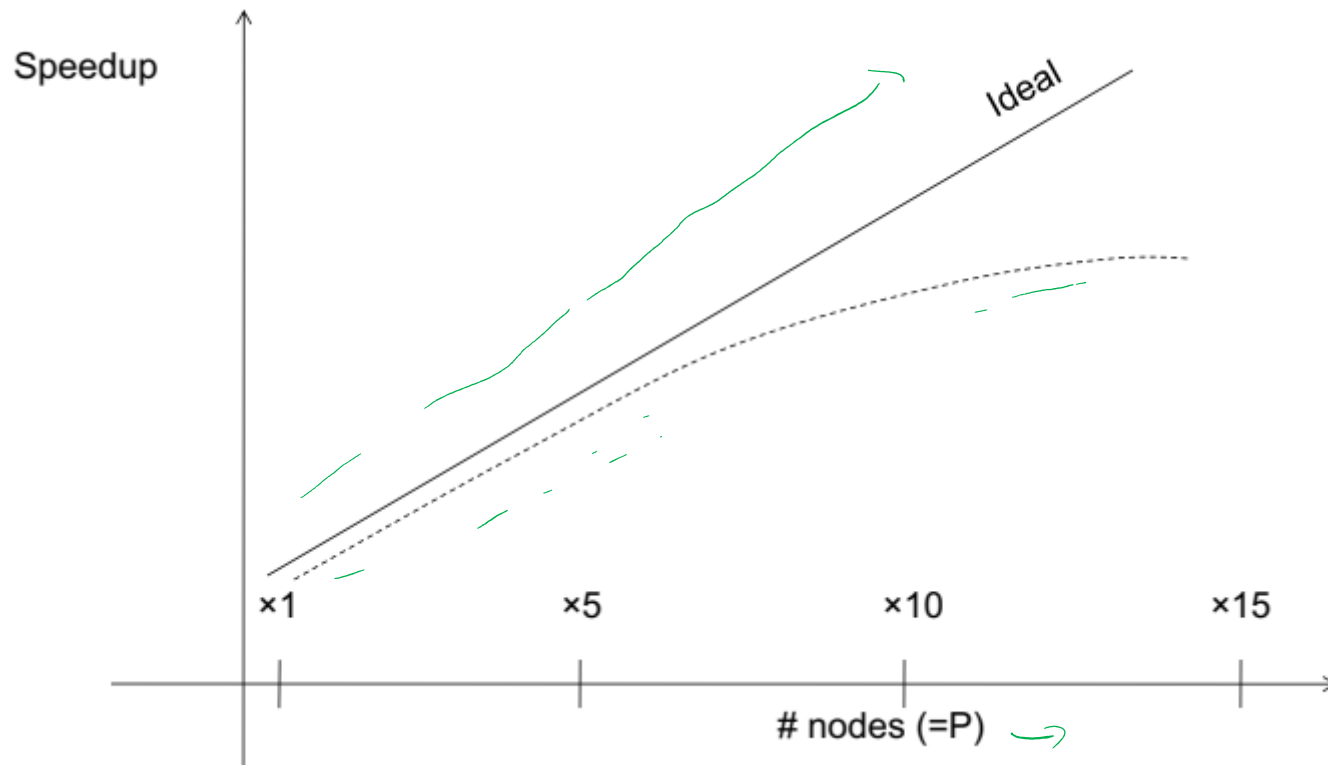
- More nodes, more data → same speed

- **OLTP:** “Speed” = transactions per second (TPS)

- **Decision Support:** “Speed” = query time



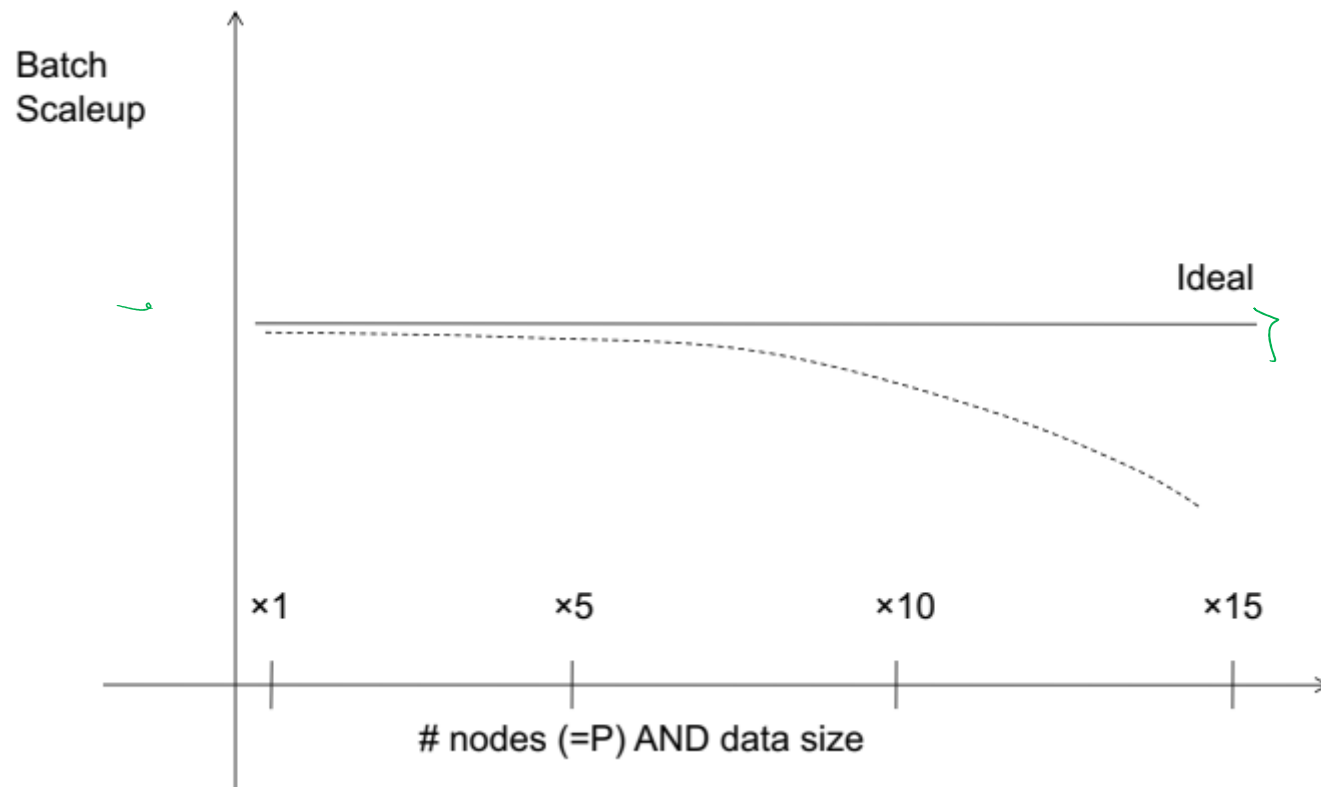
Linear vs. Non-linear Speedup



- * Speedup - More nodes, same data \rightarrow higher speed
- * Scaleup - More nodes, more data \rightarrow same speed



Linear vs. Non-linear Scaleup

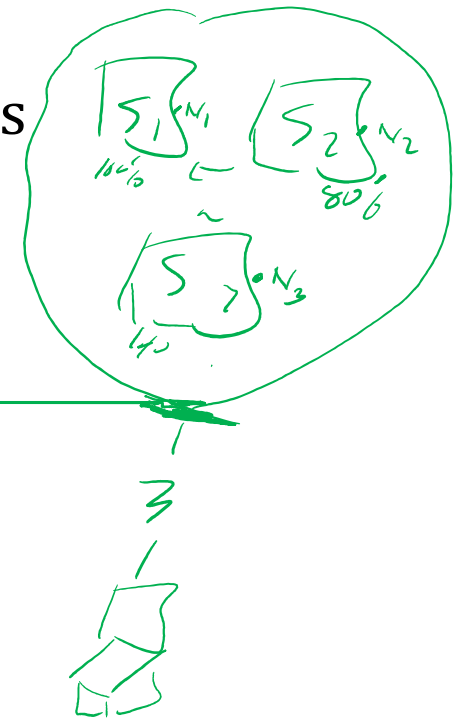


- * Speedup - More nodes, same data → higher speed
- * Scaleup - More nodes, more data → same speed



Challenges to Linear Speedup and Scaleup

- **Startup cost** –
 - Cost of starting an operation on many nodes
- **Interference** –
 - Contention for resources between nodes
- **Stragglers**
 - Slowest node becomes the bottleneck



INTRO TO DATA STRUCTURE

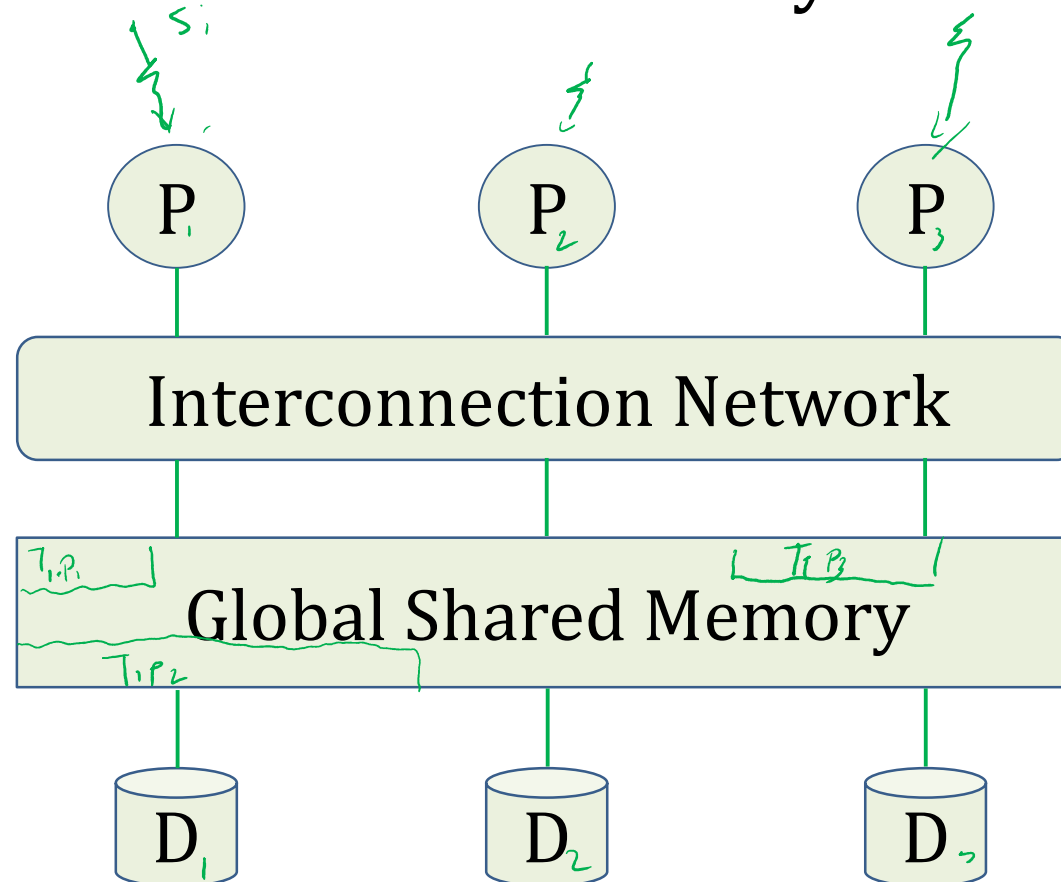
Parallel DB Architectures

Architectures for Parallel DBs

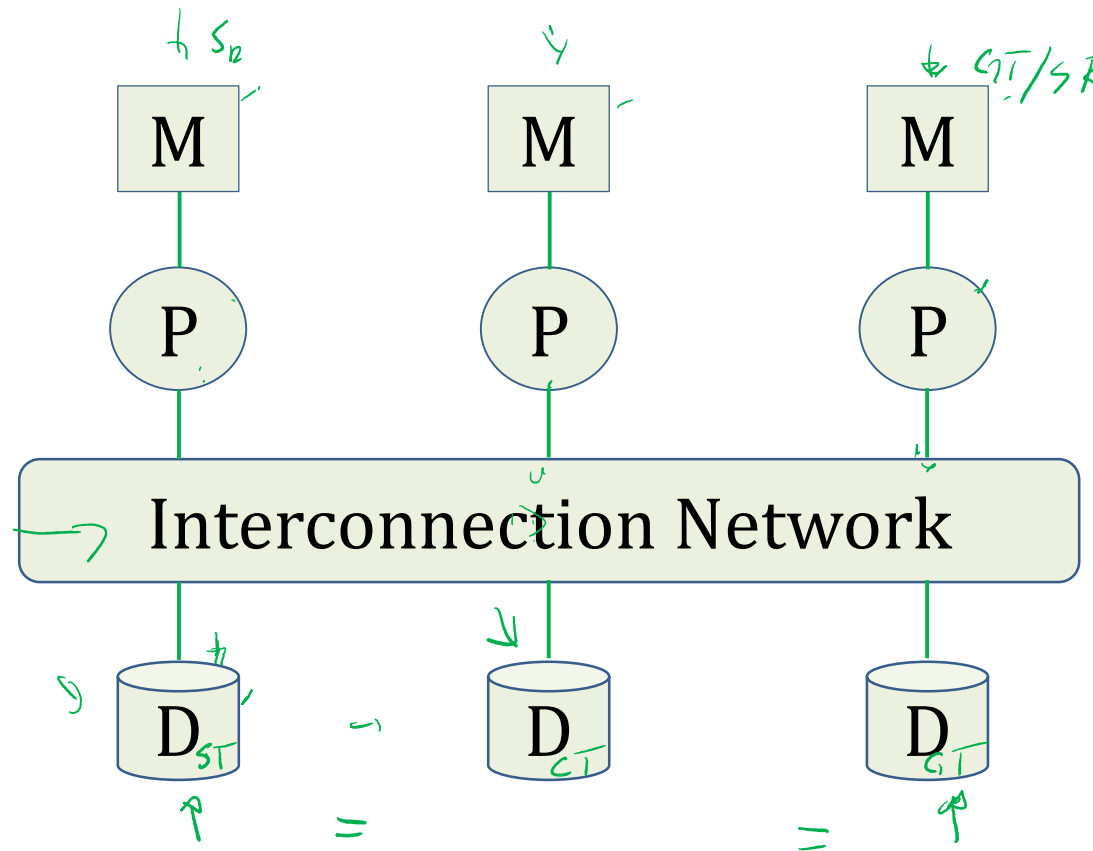
- Shared memory
- Shared disk
- Shared nothing



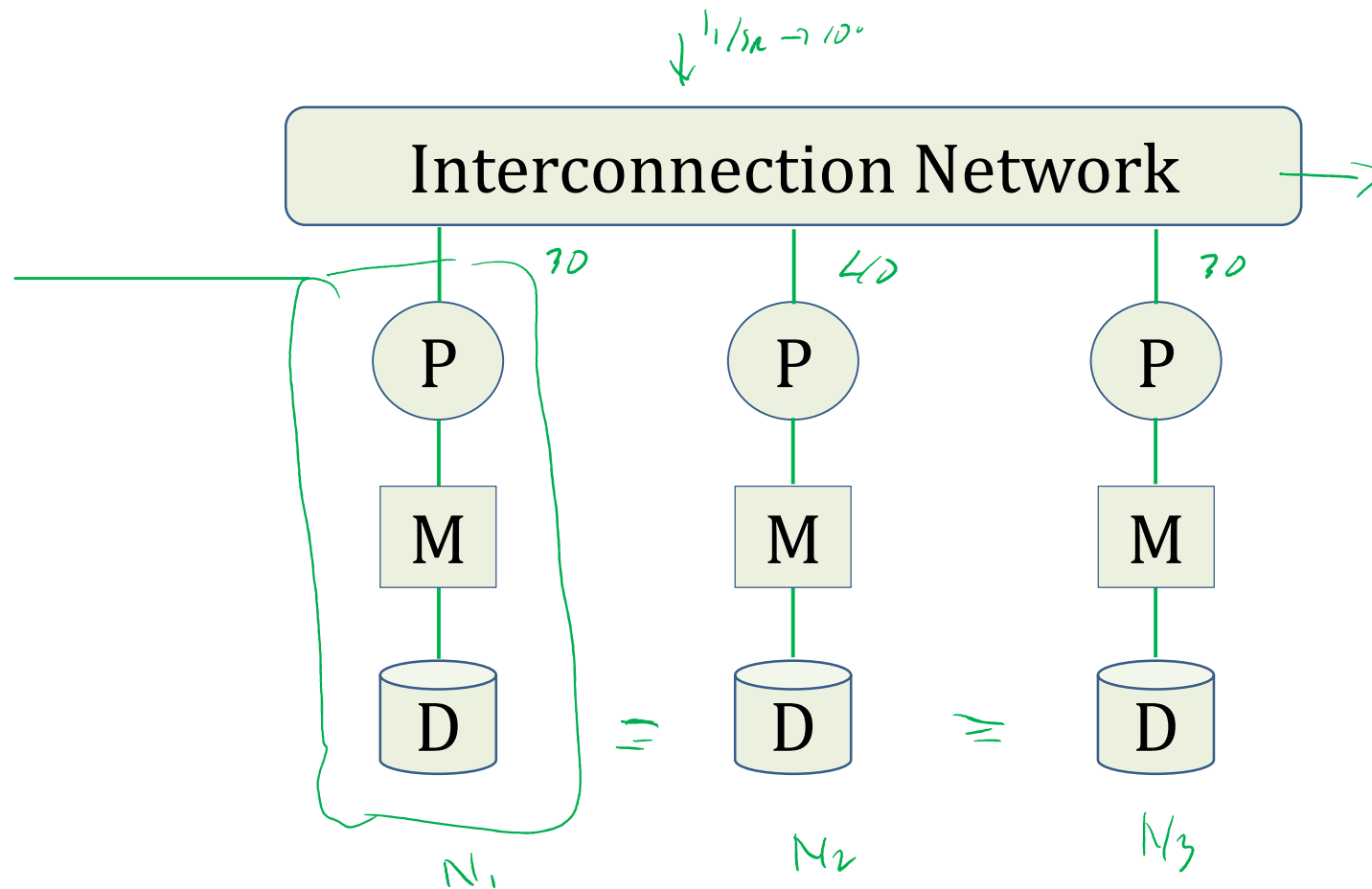
Shared Memory



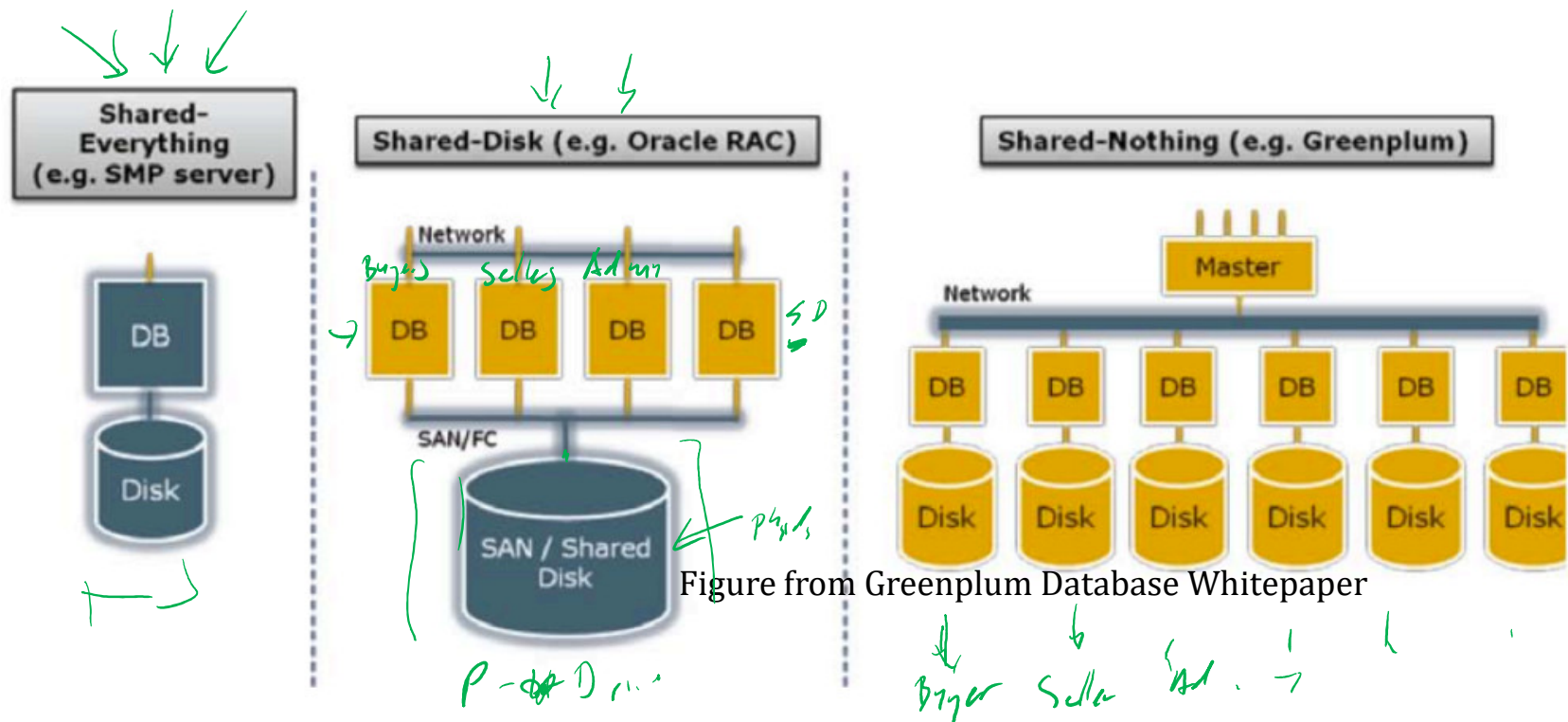
Shared Disk



Shared Nothing



A Professional View



SAN = "Storage Area Network"

Shared Memory

Random Access Memory

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

License

db server

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

- Easier to program and easy to use
- But very expensive to scale: last remaining cash cows in the hardware industry

* cash cow – provide steady income or profit that is far bigger than cost to acquire, start or upgrade (scale)



Shared Disk

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

Oracle dominates this class of systems.

Characteristics:

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



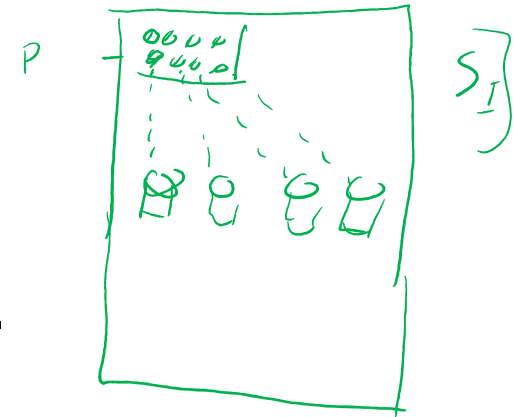
Shared Nothing

- Cluster of machines on high-speed network
- 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



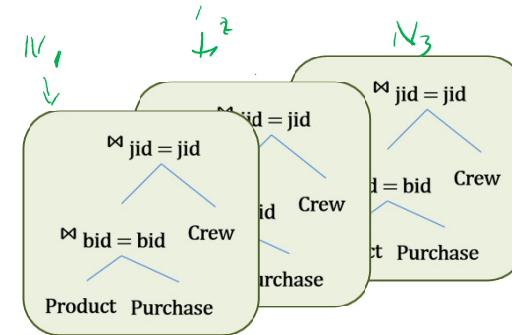
* Contention – competition for resources

INTRO TO DATA STRUCTURE

Distributed Query Processing

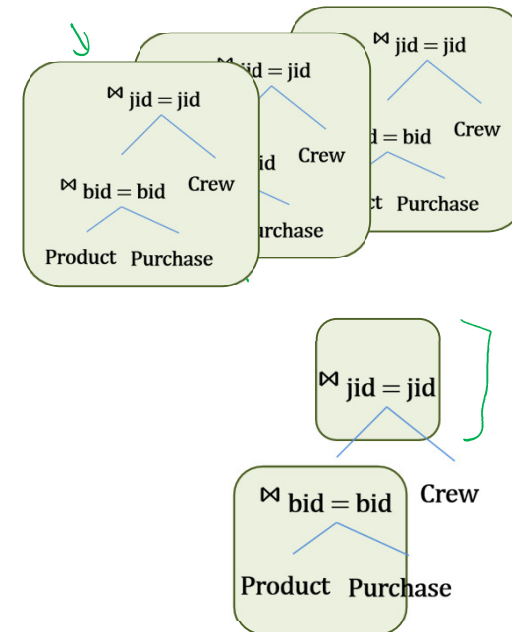
Approaches to Parallel Query Evaluation

- **Inter-query** parallelism
 - Transaction per node
 - OLTP



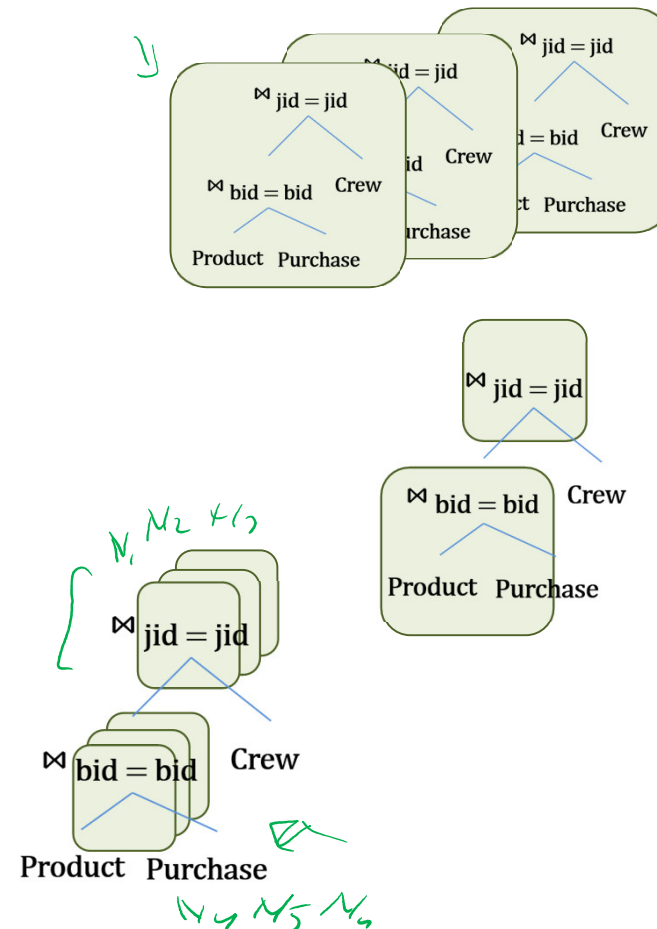
Approaches to Parallel Query Evaluation

- **Inter-query** parallelism
 - Transaction per node
 - OLTP
- **Inter-operator** parallelism
 - Operator per node
 - Both OLTP and Decision Support



Approaches to Parallel Query Evaluation

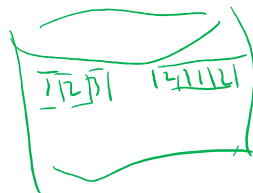
- **Inter-query** parallelism
 - Transaction per node
 - OLTP
- **Inter-operator** parallelism
 - Operator per node
 - Both OLTP and Decision Support
- **Intra-operator** parallelism
 - Operator on multiple node
 - Decision Support
 - Most scalable



Review on Single Node Query Processing

Given relations R(A, B) and S(B, C), no indexes:

- **Selection:** $\sigma_{A=123}(R)$
 - Scan file R, select records with A=123
- **Group-by:** $\gamma_{A, \text{sum}(B)}(R)$
 - Scan file R, insert into a hash table using attribute A as key
 - When a new key is equal to an existing one, add B to value
- **Join:** $R \bowtie S$
 - Scan file S, insert into a hash table using attribute B as key
 - Scan file R, probe the hash table using attribute B

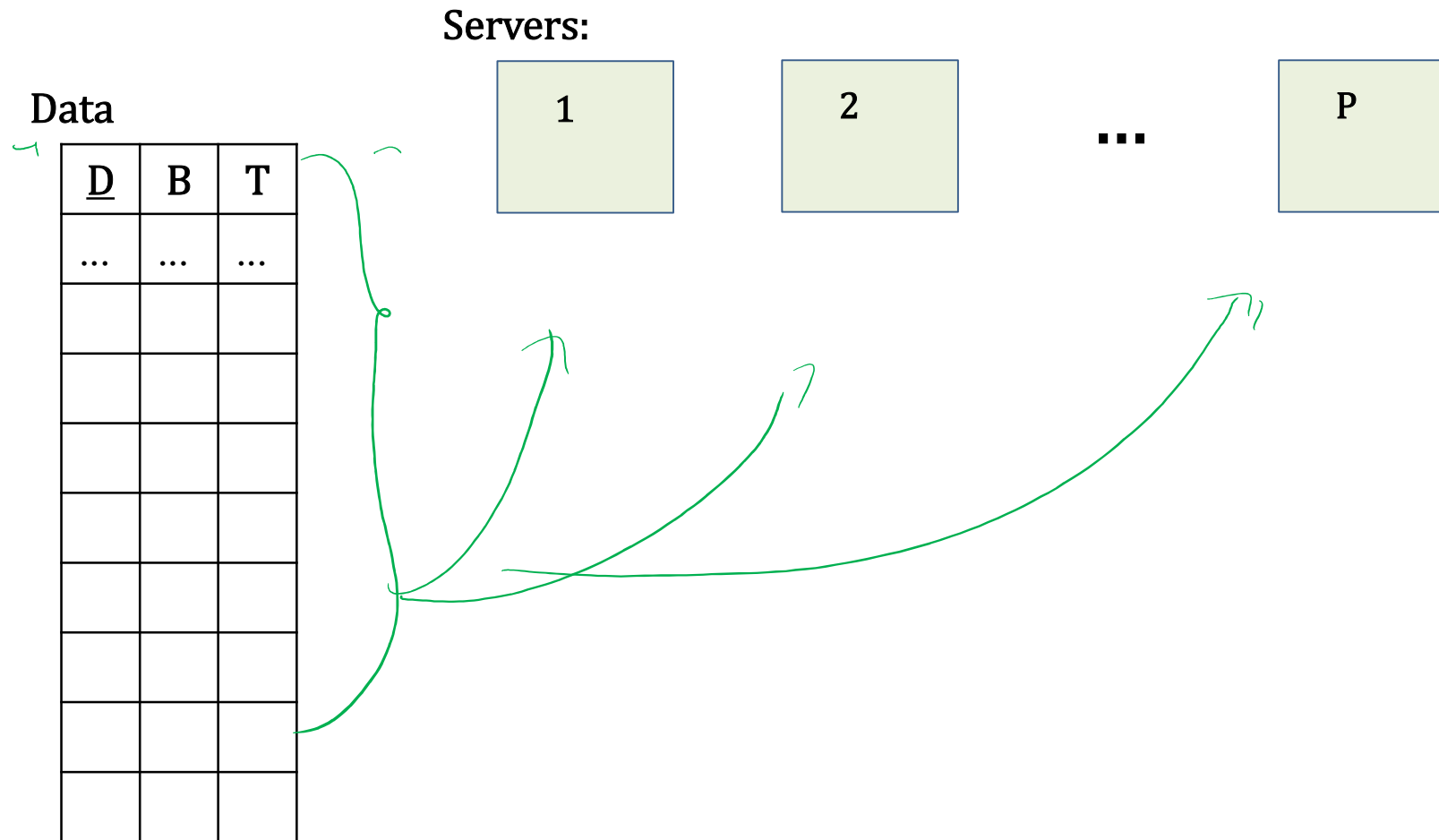


Distributed Query Processing

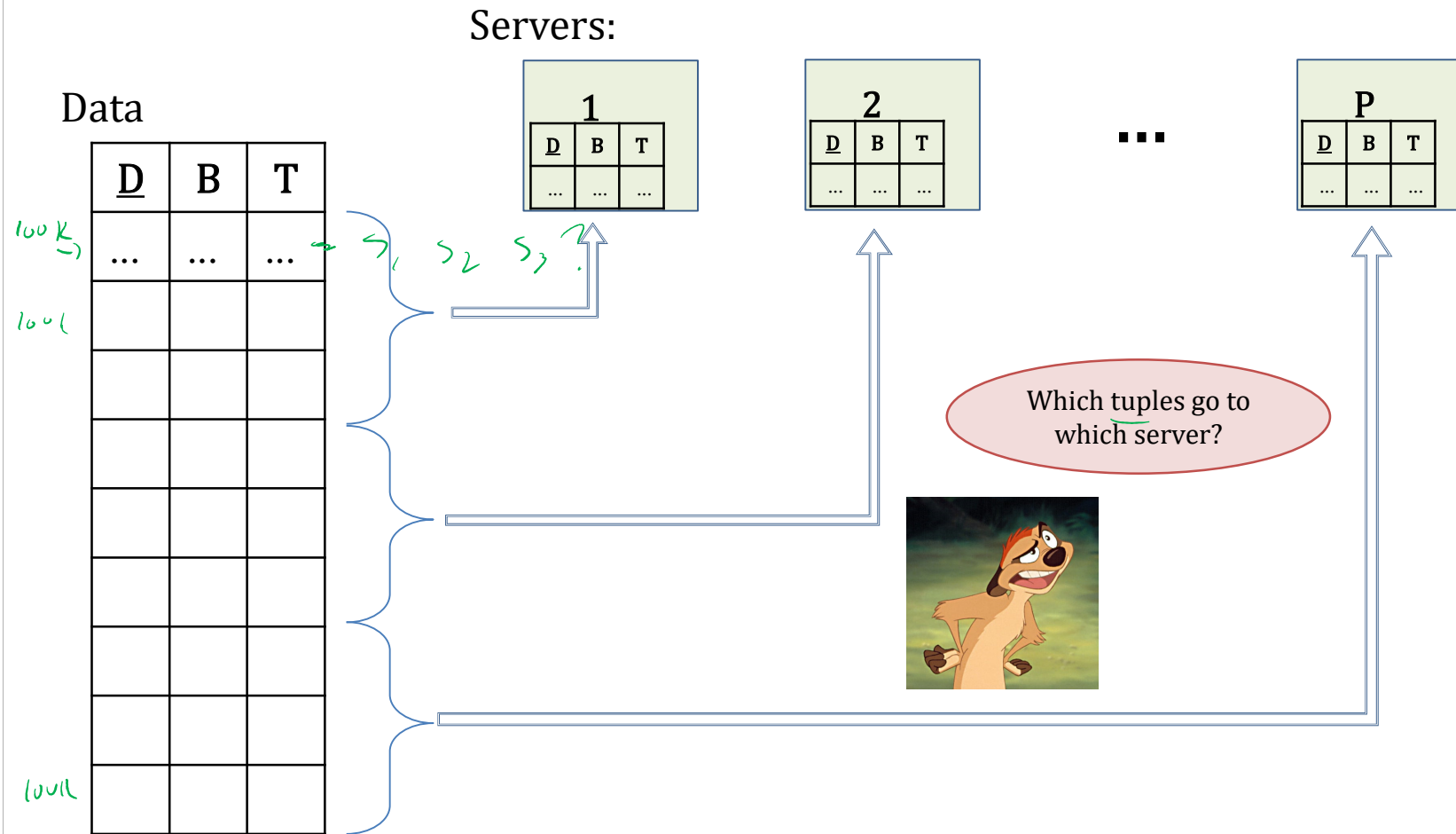
- Data is horizontally partitioned across many servers
- Operators may require data reshuffling
 - Not all the needed data is in one place



Horizontal Data Partitioning



Horizontal Data Partitioning



Horizontal Data Partitioning

- Block Partition:

- Partition tuples arbitrarily s.t. $\text{size}(R_1) \approx \dots \approx \text{size}(R_p)$



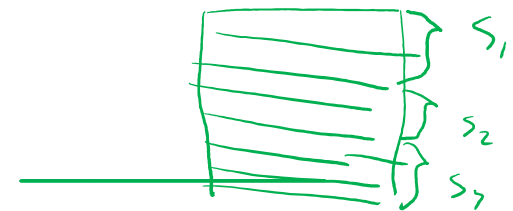
- Hash partitioned on attribute A:

- Tuple t goes to chunk i , where $i = h(t.A) \bmod P + 1$

$$P \bmod S$$

- Range partitioned on attribute A:

- Partition the range of A into $-\infty = v_0 < v_1 < \dots < v_p = \infty$
- Tuple t goes to chunk i , if $v_{i-1} < t.A < v_i$



Parallel GroupBy

Data: $R(\underline{K}, A, B, C)$

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

- How can we compute in each case
 - R is hash-partitioned on A
 - R is block-partitioned
 - R is hash-partitioned on K

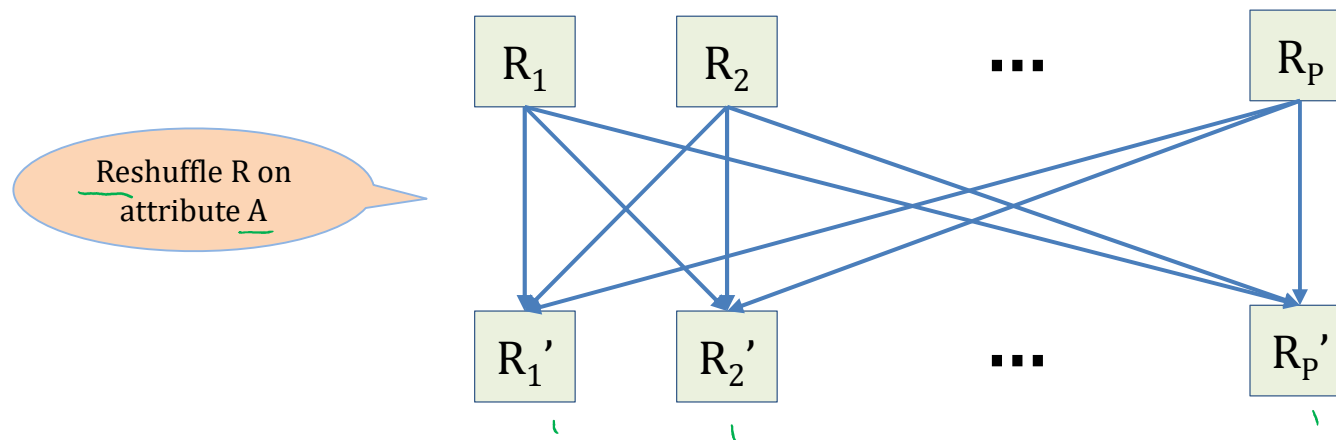


Parallel GroupBy

Data: $R(\underline{K}, A, B, C)$

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

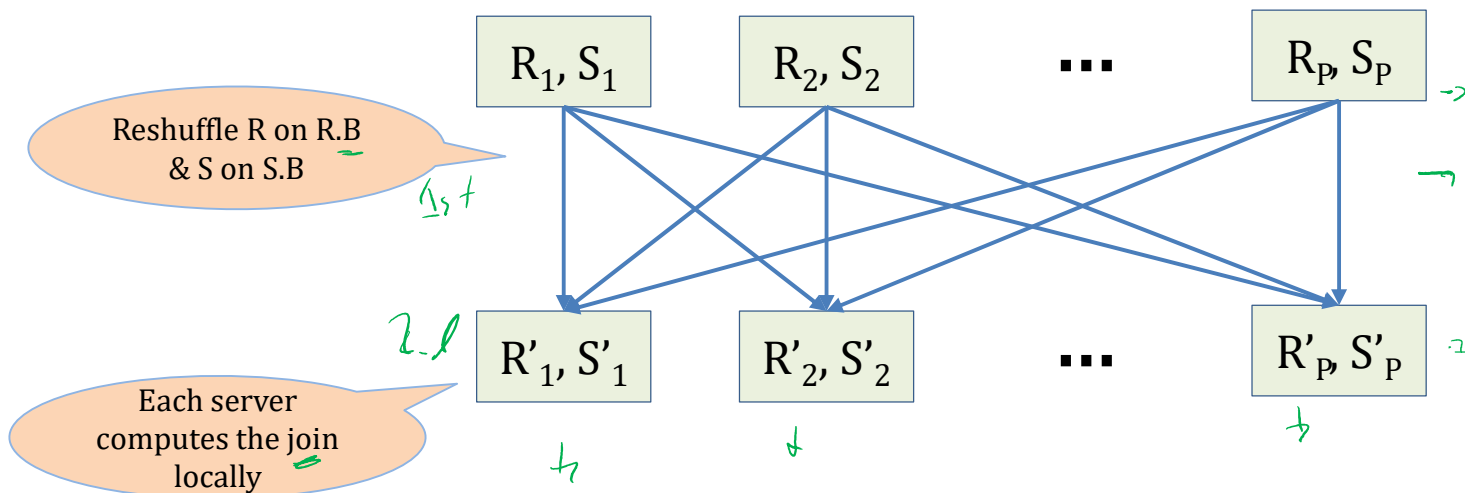
- R is block-partitioned or hash-partitioned on K



Parallel Join

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

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



Data: $R(K1, A, B), S(K2, B, C)$

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

Parallel Join



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

Parallel Join

Range
Partition

R1		S1	
K1	B	K2	B
1	20	101	50
2	50	102	50

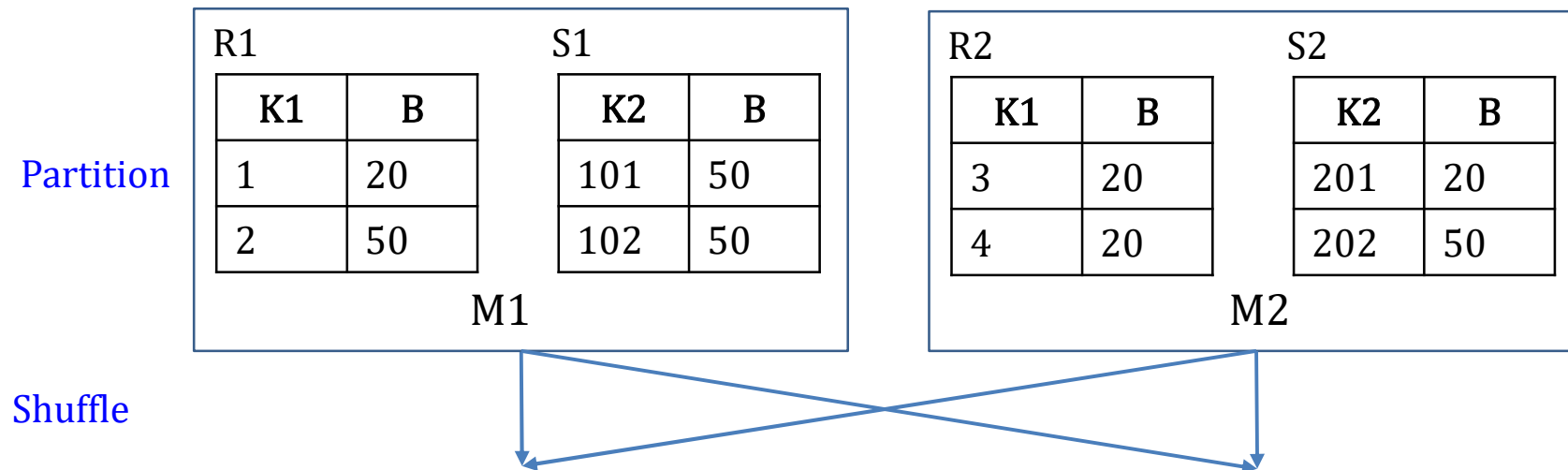
M1

R2		S2	
K1	B	K2	B
3	20	201	20
4	20	202	50

M2

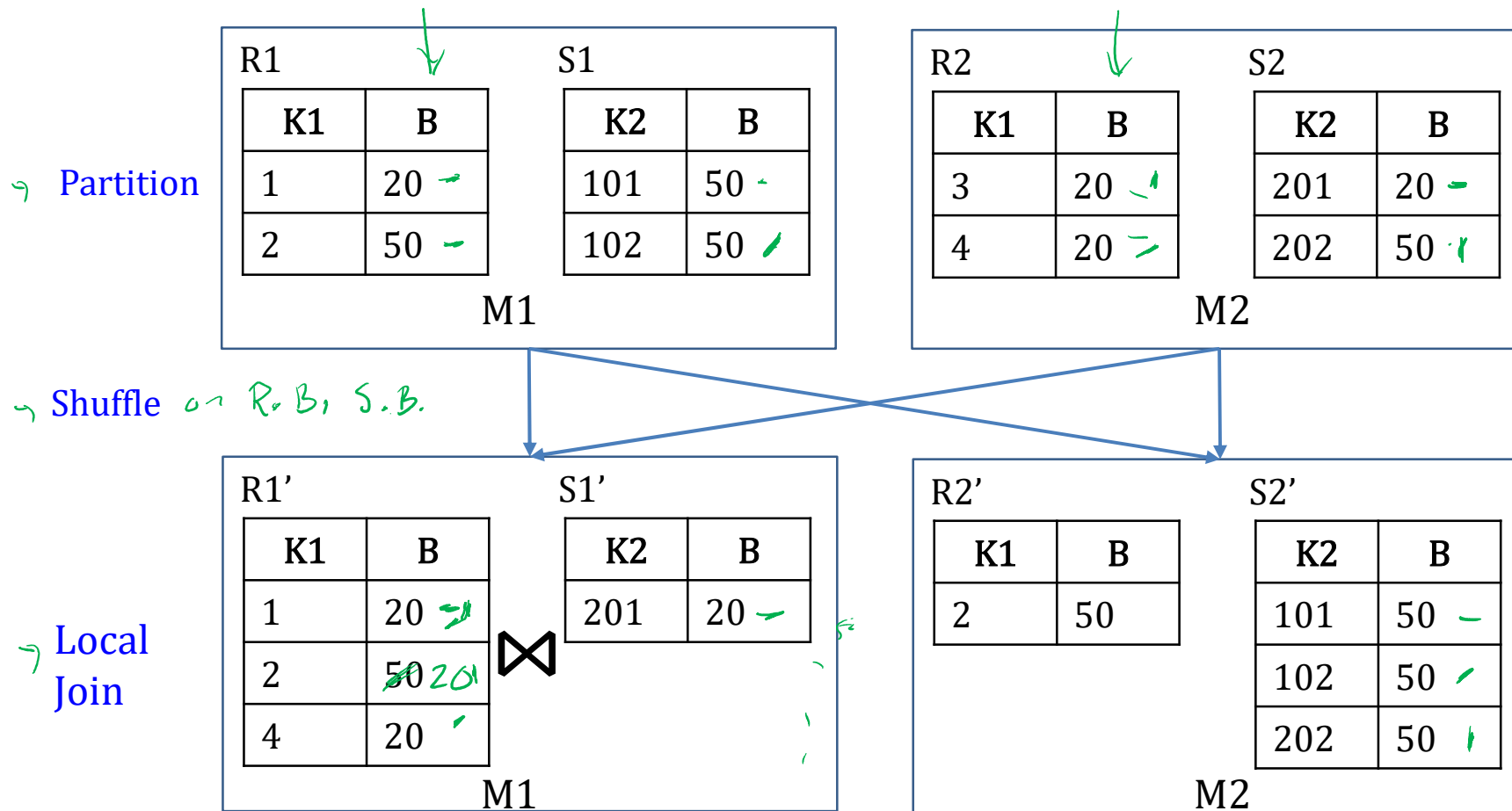
Data: $R(K1, A, B), S(K2, B, C)$
Query: $R(K1, A, B) \bowtie S(K2, B, C)$

Parallel Join



Data: R(K1, A, B), S(K2, B, C)
Query: R(K1, A, B) ⋈ S(K2, B, C)

Parallel Join



Speedup and Scaleup

- Consider:
 - Query: $\gamma_{A, \text{sum}}(C)(R) \rightarrow \gamma_{A, \text{sum}(C)}(R)$
 - Runtime: dominated by reading chunks from disk

Speedup and Scaleup

- Consider:
 - Query: $\gamma_{A, \text{sum}}(C)(R)$
 - 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 $\frac{1}{2}$ as many chunks)



Speedup and Scaleup

• Consider:

- Query: $\gamma_{A, \text{sum}}(C)(R)$
- 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 $\frac{1}{2}$ 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)

as double number of nodes \Rightarrow same running time (t_1)
 B' same running time (t_1)
 C' half running time $(\frac{t_1}{2})$



Uniform Data vs. Skewed Data

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

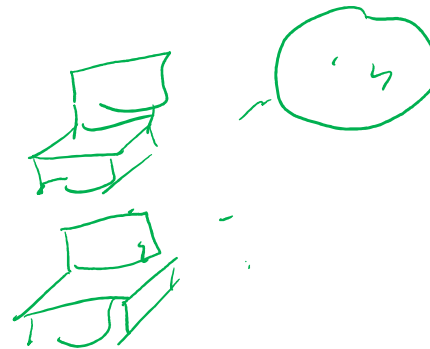
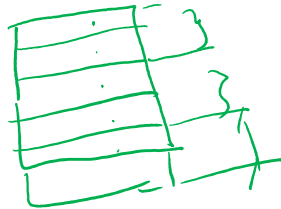
• Block partition

• Hash-partition

- On the key K
- On the attribute A

• Range-partition

- On the key K
- On the attribute A



$PID \% 5$



Uniform Data vs. 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

Uniform

Uniform

May be skewed

May be skewed

Assuming good hash function

Ex.: When all tuples have the same value of the attribute A, then all records end up in the same partition

Difficult to partition the range of A uniformly

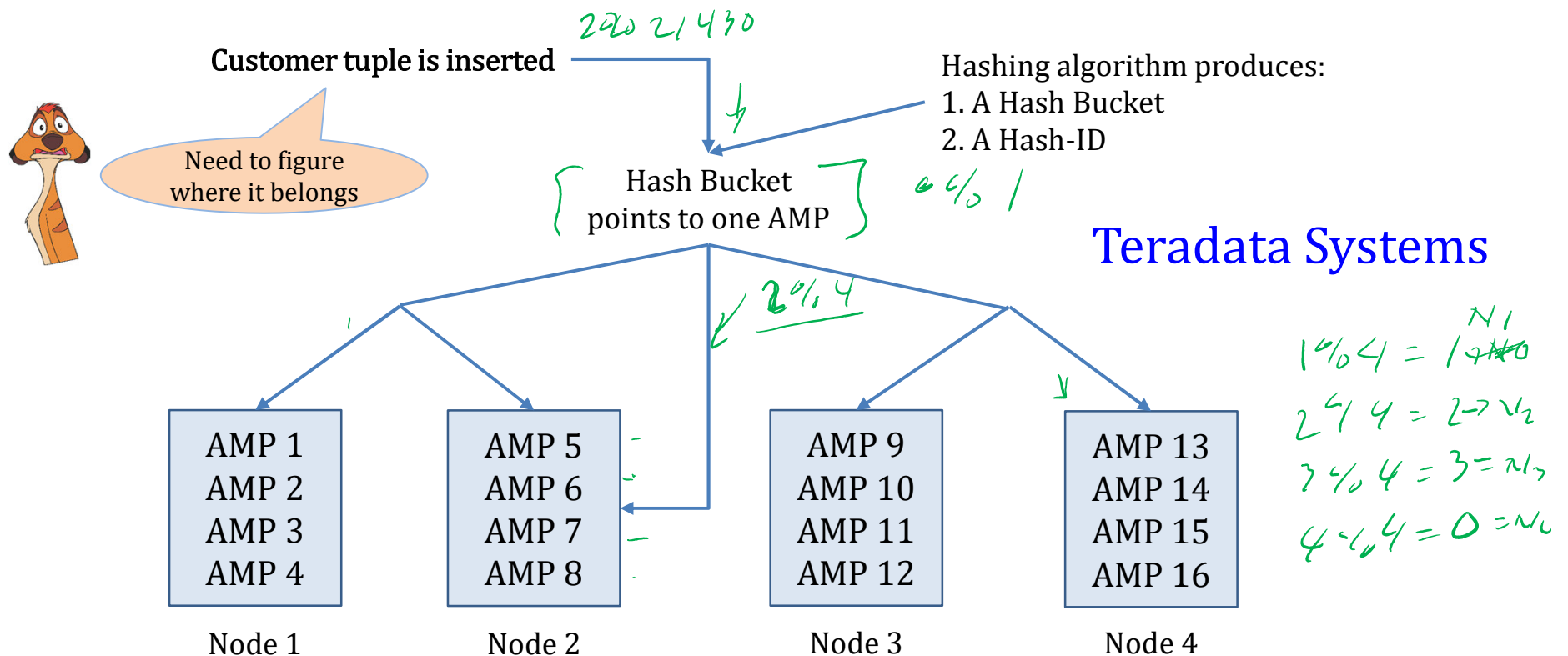
Loading Data into a Parallel DBMS

Customer tuple is inserted



Need to figure
where it belongs

Loading Data into a Parallel DBMS

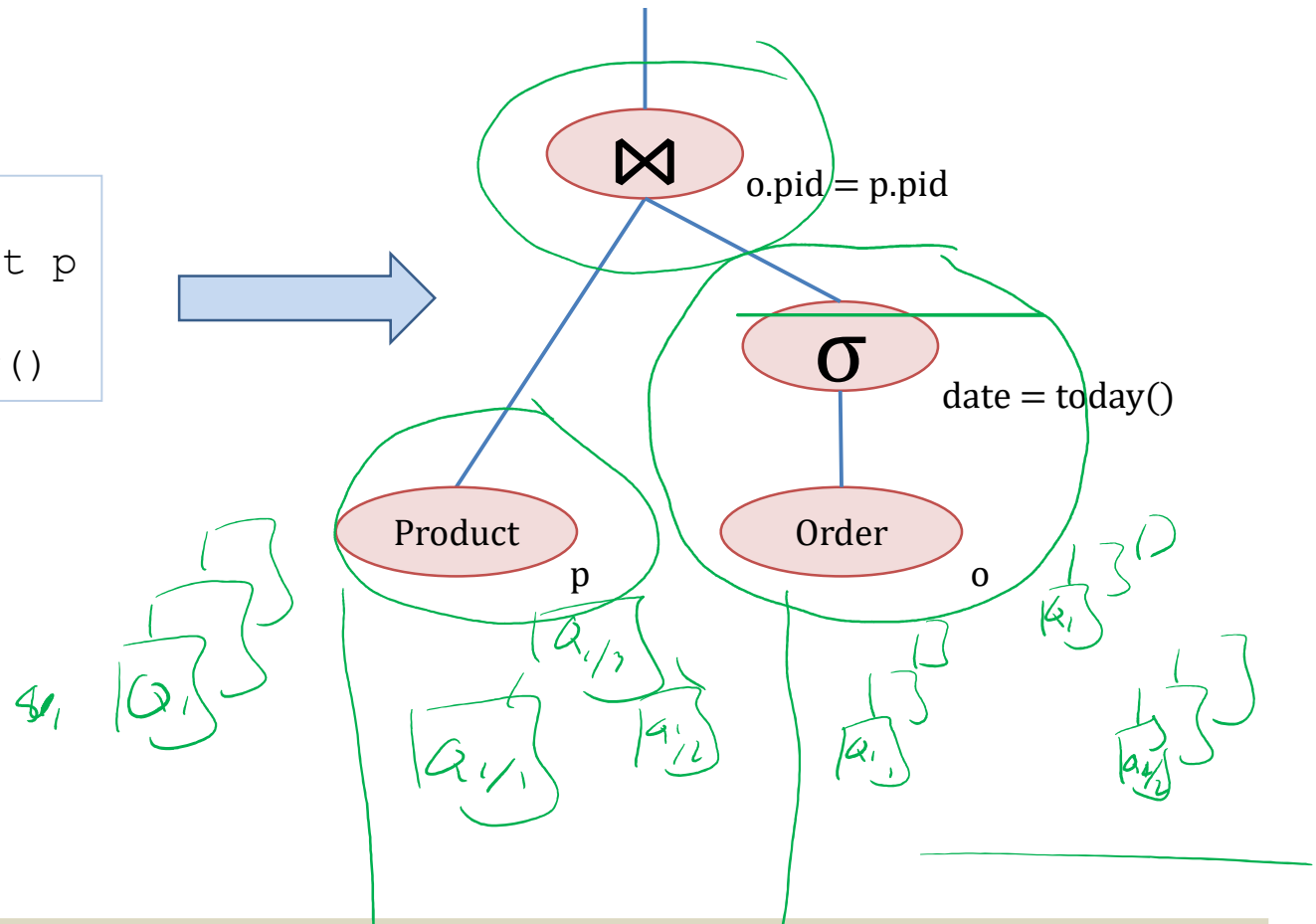


AMP = "Access Module Processor" = unit of parallelism

Example Parallel Query Execution

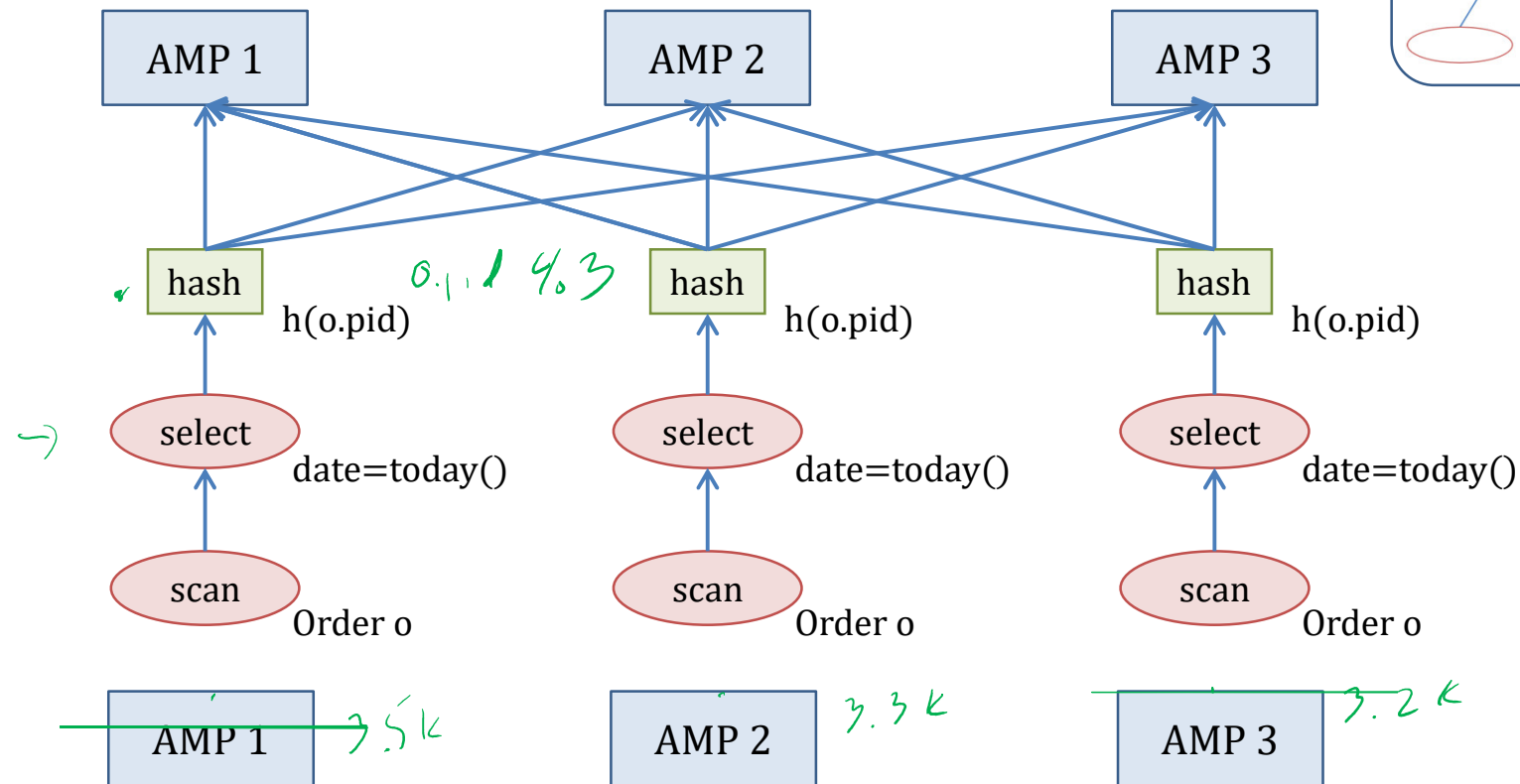
Order(oid, pid, date)
Product(pid, ...)

```
SELECT *
FROM Order o, Product p
WHERE o.pid = p.pid
      AND o.date = today()
```



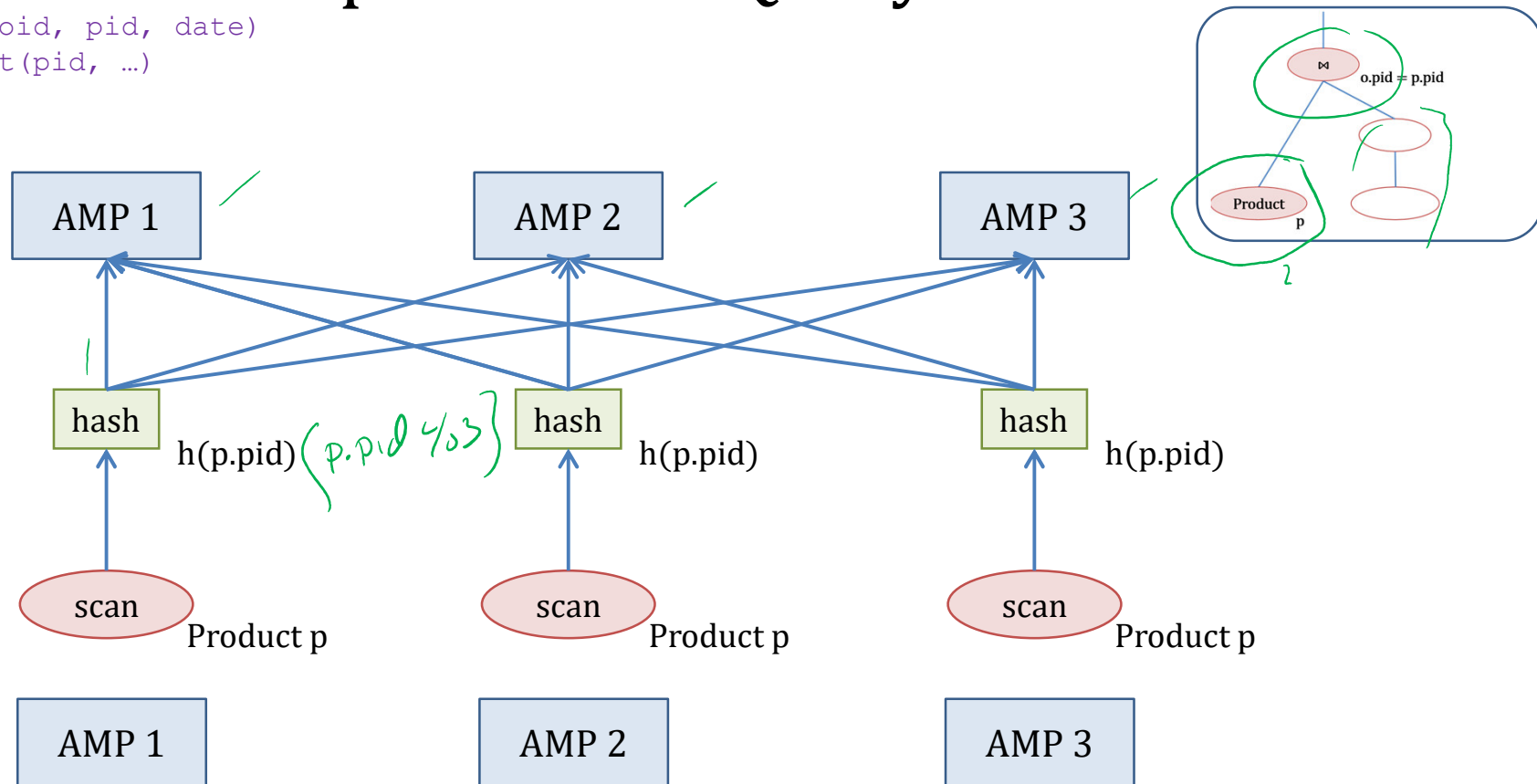
Example Parallel Query Execution

Order(oid, pid, date)
Product(pid, ...)



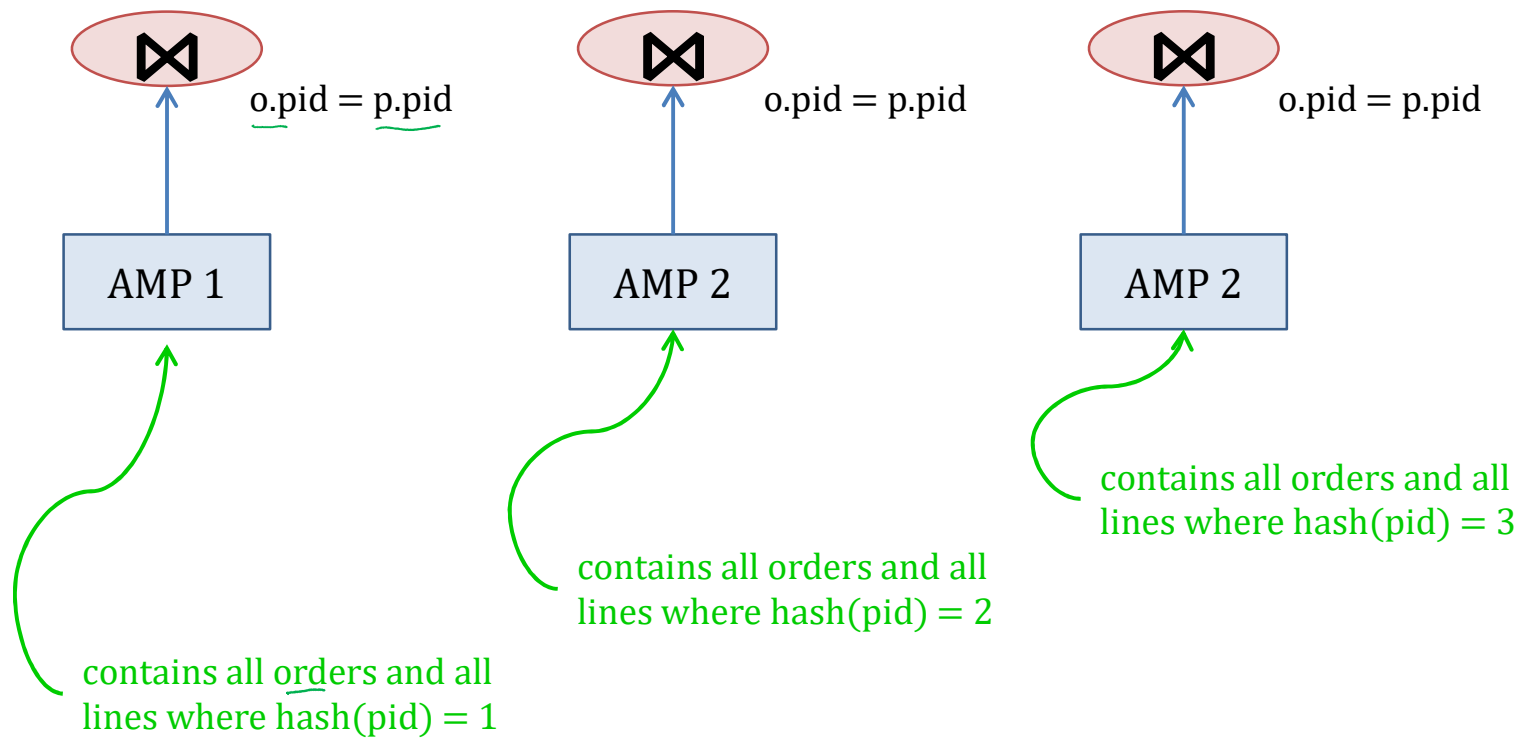
Example Parallel Query Execution

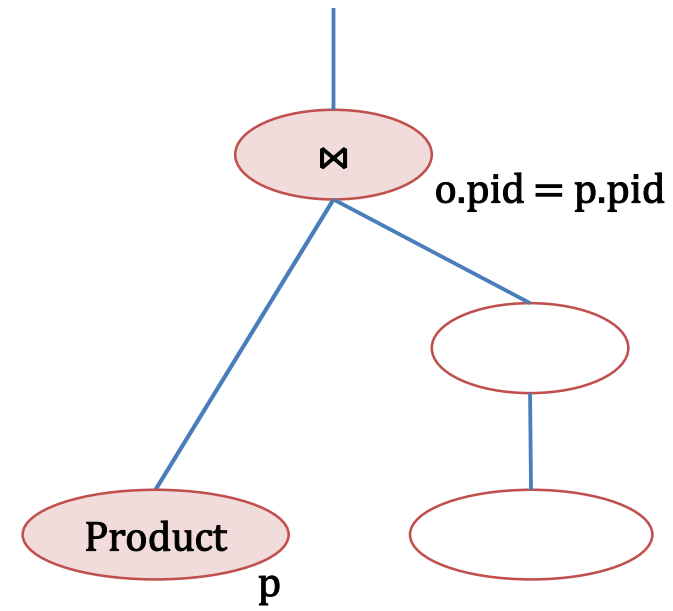
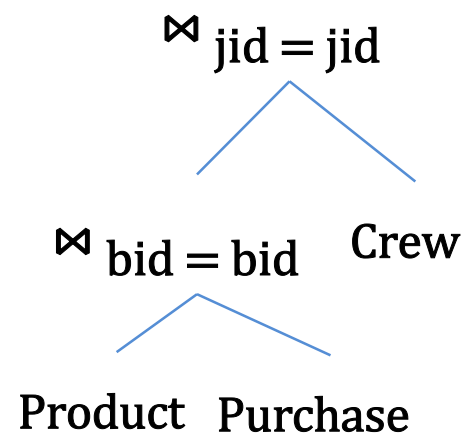
Order(oid, pid, date)
Product(pid, ...)



Example Parallel Query Execution

```
Order(oid, pid, date)  
Product(pid, ...)
```





Thank you.