# ADMINISTRIVIA

**Monday Dec 4$^{th}$ – NuoDB**
→ Barry Morris (Co-Founder, Exec. Chairman)

**Wednesday Dec 6$^{th}$ – Potpourri + Final Review**
→ Vote for what system you want me to talk about.
→ http://cmudb.io/f17-systems

**Wednesday Dec 6$^{th}$ – Project #4**

# OLTP VS. OLAP

**On-line Transaction Processing (OLTP):**
→ Short-lived txns.
→ Small footprint.
→ Repetitive operations.

**On-line Analytical Processing (OLAP):**
→ Long running queries.
→ Complex joins.
→ Exploratory queries.

CARNEGIE MELLON
**DATABASE GROUP**

# TODAY'S AGENDA

Partitioning

Distributed Join Algorithms

CARNEGIE MELLON
**DATABASE GROUP**

# DATABASE PARTITIONING

Split database across multiple resources:
→ Disks, nodes, processors.
→ Sometimes called "sharding"

The DBMS executes query fragments on each partition and then combines the results to produce a single answer.

# NAÏVE TABLE PARTITIONING

Each node stores one and only table.

Assumes that each node has enough storage space for a table.

CARNEGIE MELLON
DATABASE GROUP

# NAÏVE TABLE PARTITIONING

Table1                Table2                Partitions

**Ideal Query:**

```
SELECT * FROM table
```

CARNEGIE MELLON
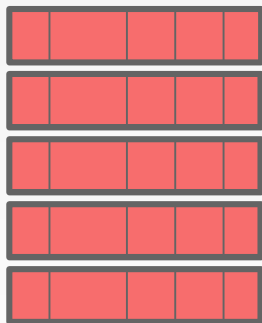**DATABASE GROUP**
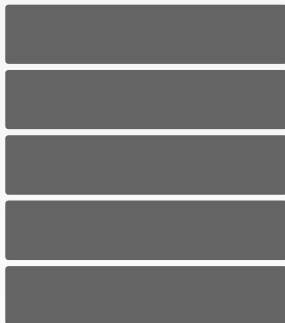
# NAÏVE TABLE PARTITIONING

Table1
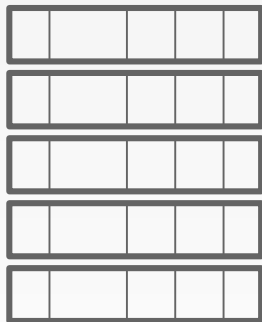
Table2

Partitions

**Ideal Query:**

```
SELECT * FROM table
```
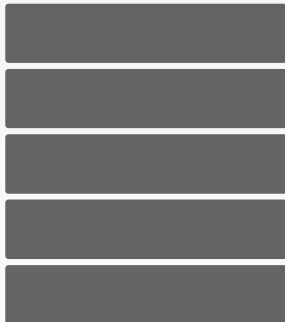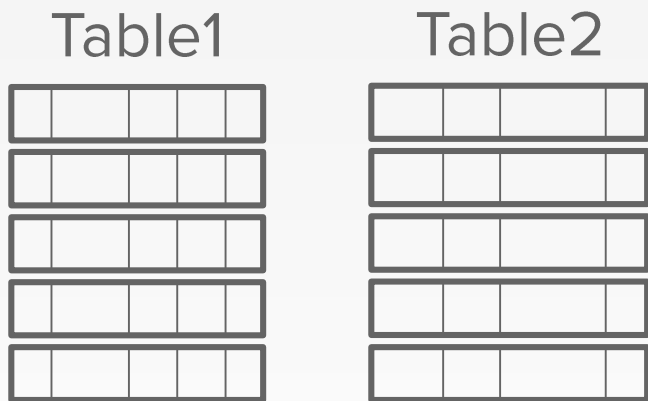
# NAÏVE TABLE PARTITIONING

Table1

Table2

Partitions

**Ideal Query:**

```
SELECT * FROM table
```

Table1

CARNEGIE MELLON
DATABASE GROUP

# NAÏVE TABLE PARTITIONING

Table1

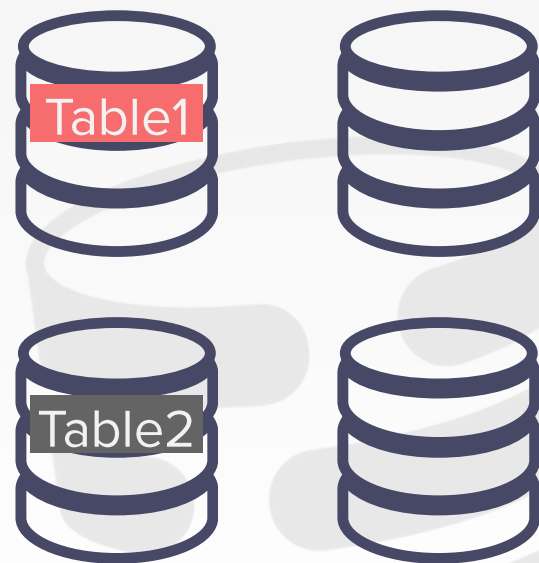Table2

Partitions



**Ideal Query:**

```
SELECT * FROM table
```

CARNEGIE MELLON
DATABASE GROUP

# HORIZONTAL PARTITIONING

Split a table's tuples into disjoint subsets.
→ Choose column(s) that divides the database equally in terms of size, load, or usage.
→ Each tuple contains all of its columns.
→ Hash Partitioning, Range Partitioning

The DBMS can partition a database **physical** (shared nothing) or **logically** (shared disk).

CARNEGIE MELLON
DATABASE GROUP

# HORIZONTAL PARTITIONING

Partitioning Key

Table1

| 101 | a | XXX | 2017-11-29 |
| 102 | b | XXY | 2017-11-28 |
| 103 | c | XYZ | 2017-11-29 |
| 104 | d | XYX | 2017-11-27 |
| 105 | e | XYY | 2017-11-29 |

Partitions

**Ideal Query:**

```
SELECT * FROM table
 WHERE partitionKey = ?
```

CARNEGIE MELLON
DATABASE GROUP

# HORIZONTAL PARTITIONING

Partitioning Key
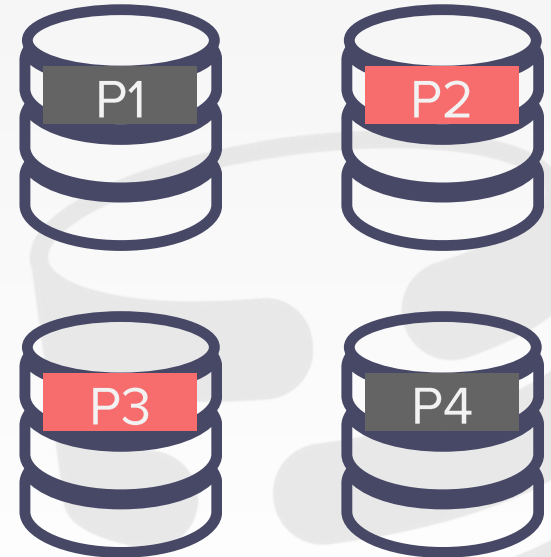
Table1

101 a XXX 2017-11-29    hash(a)%4 = P2

102 b XXY 2017-11-28    hash(b)%4 = P4

103 c XYZ 2017-11-29    hash(c)%4 = P3

104 d XYX 2017-11-27    hash(d)%4 = P2

105 e XYY 2017-11-29    hash(e)%4 = P1

Partitions

**Ideal Query:**

```
SELECT * FROM table
 WHERE partitionKey = ?
```

CARNEGIE MELLON
DATABASE GROUP

# HORIZONTAL PARTITIONING

Partitioning Key

Table1

| 101 | a | XXX | 2017-11-29 |
| 102 | b | XXY | 2017-11-28 |
| 103 | c | XYZ | 2017-11-29 |
| 104 | d | XYX | 2017-11-27 |
| 105 | e | XYY | 2017-11-29 |

hash(a)%4 = P2

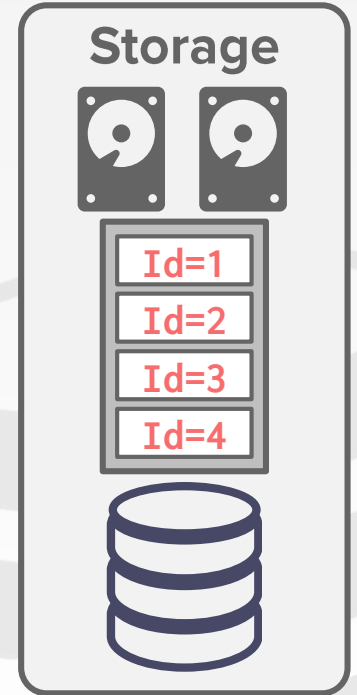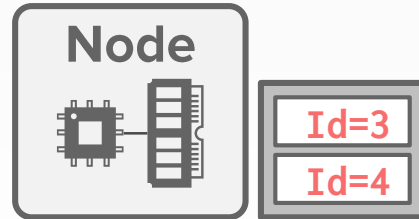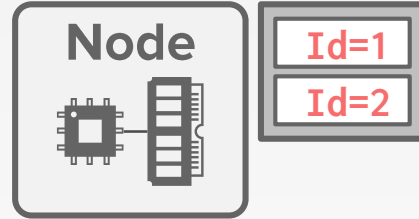hash(b)%4 = P4

hash(c)%4 = P3

hash(d)%4 = P2

hash(e)%4 = P1
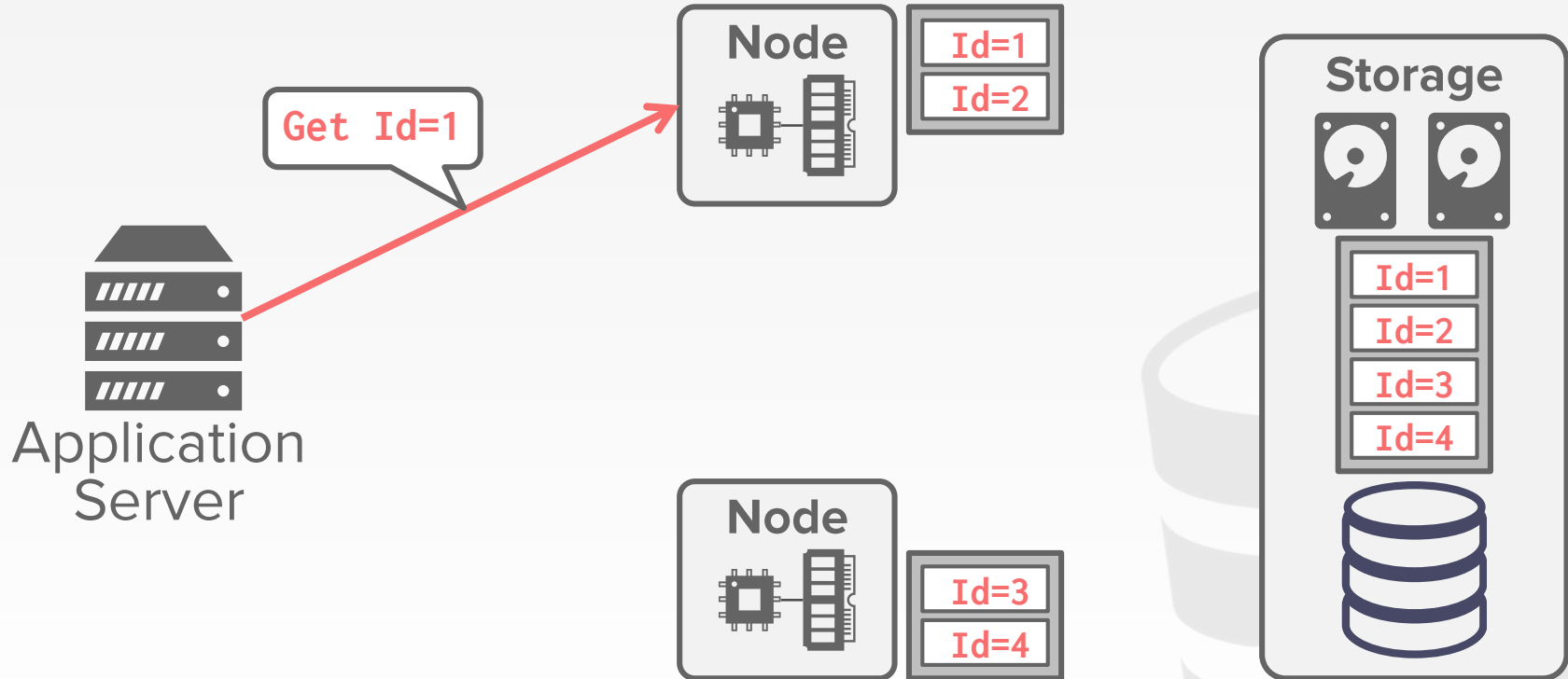
**Ideal Query:**

```
SELECT * FROM table
 WHERE partitionKey = ?
```
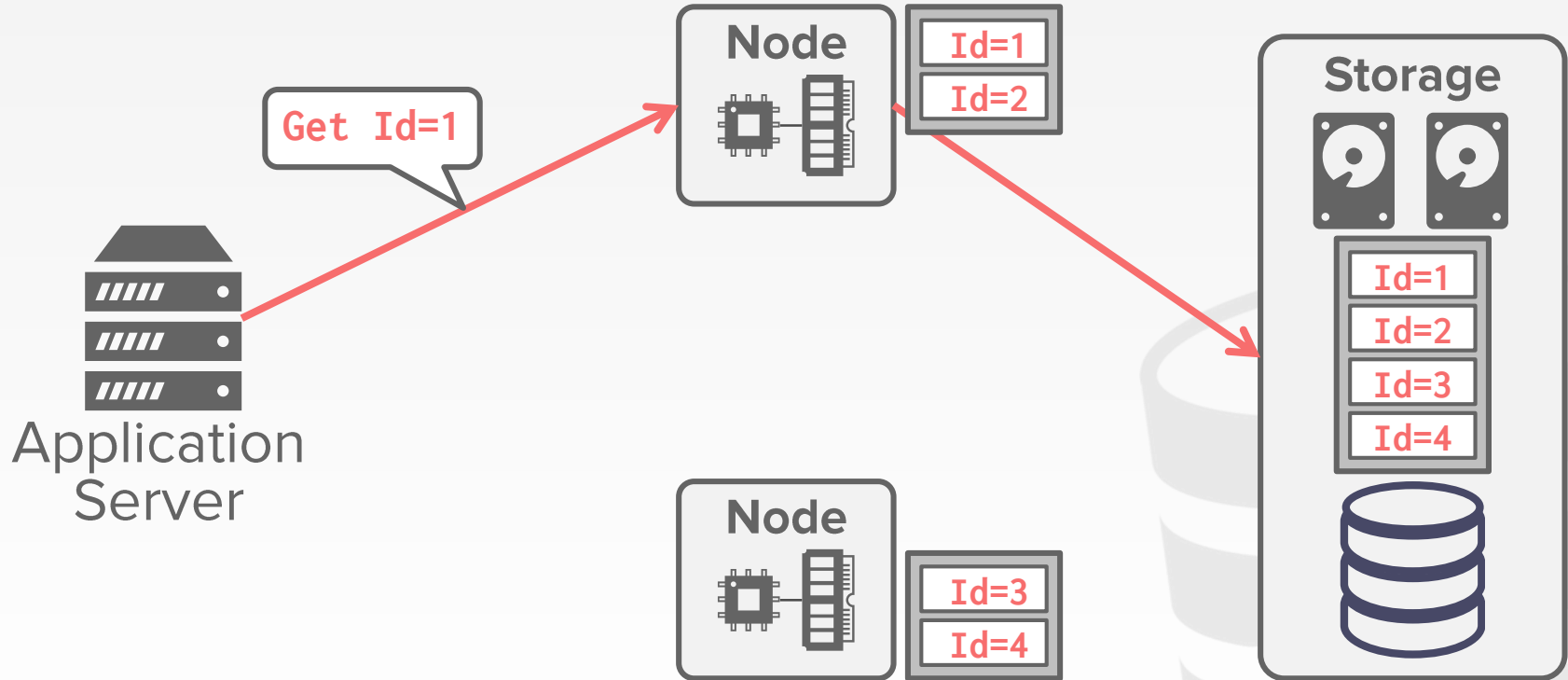
Partitions

P1    P2

P3    P4

CARNEGIE MELLON
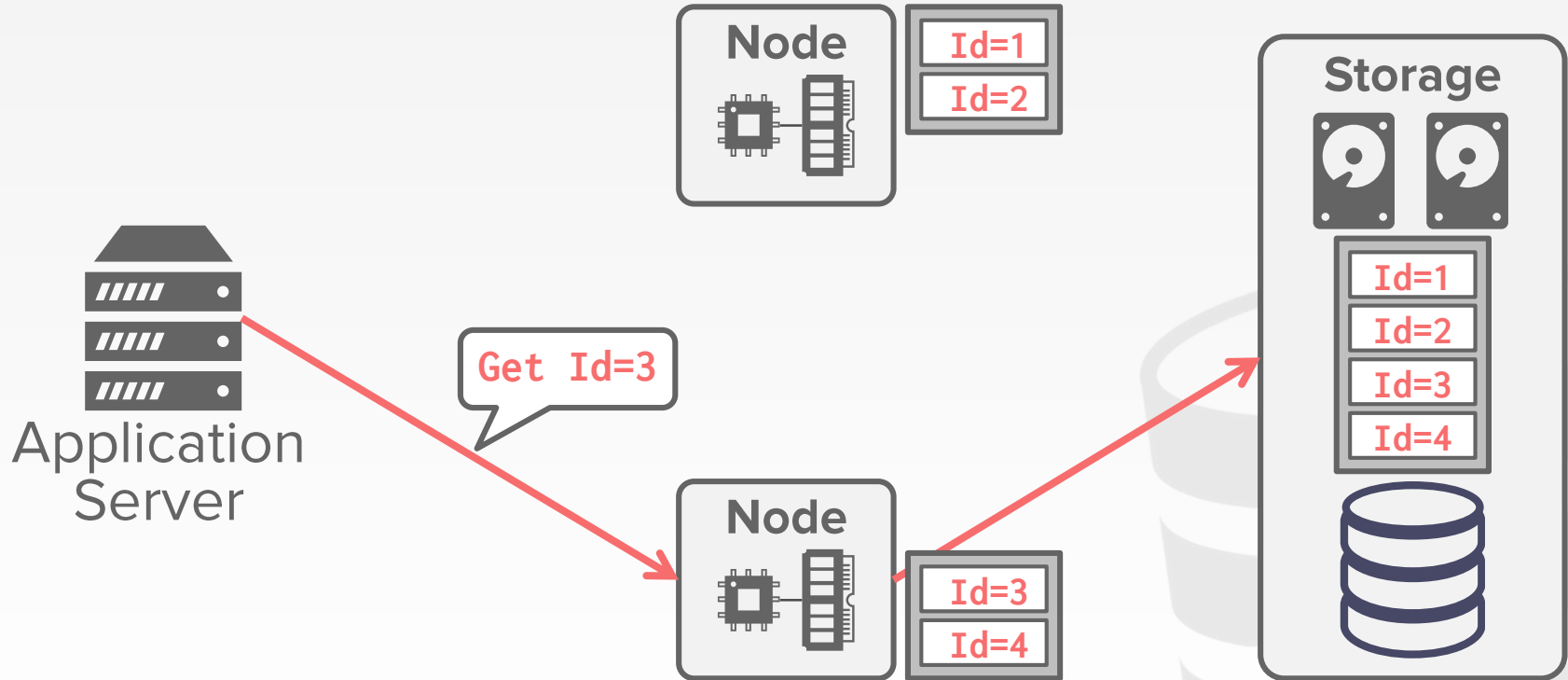DATABASE GROUP

# LOGICAL PARTITIONING

# LOGICAL PARTITIONING

# LOGICAL PARTITIONING

# LOGICAL PARTITIONING

# PHYSICAL PARTITIONING
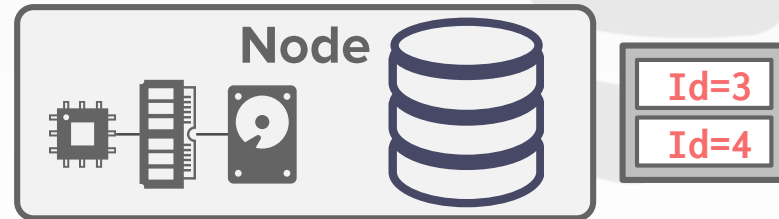
# PHYSICAL PARTITIONING

# OBSERVATION

The efficiency of a distributed join depends on the target tables' partitioning schemes.

One approach is to put entire tables on a single node and then perform the join.
→ You lose the parallelism of a distributed DBMS.
→ Costly data transfer over the network.

CARNEGIE MELLON
DATABASE GROUP

# DISTRIBUTED JOIN ALGORITHMS

To join tables **A** and **B**, the DBMS needs to get the proper tuples on the same node.

Once there, it then executes the same join algorithms that we discussed earlier in the semester.

# SCENARIO #1

One table is replicated at every node.
Each node joins its local data and then
sends their results to a coordinating node.

```
SELECT * FROM T1, T2
WHERE T1.id = T2.id
```

Id:1-100     T1 (Id)          T1 (Id)     Id:101-200

*Replicated*     T2          T2     *Replicated*

CARNEGIE MELLON
DATABASE GROUP

# SCENARIO #1

One table is replicated at every node. Each node joins its local data and then sends their results to a coordinating node.

```
SELECT * FROM T1, T2
WHERE T1.id = T2.id
```

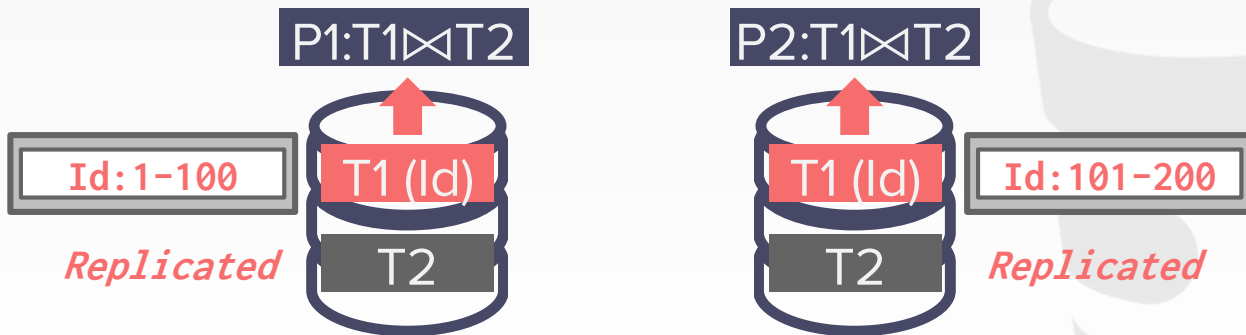# SCENARIO #1

One table is replicated at every node. Each node joins its local data and then sends their results to a coordinating node.
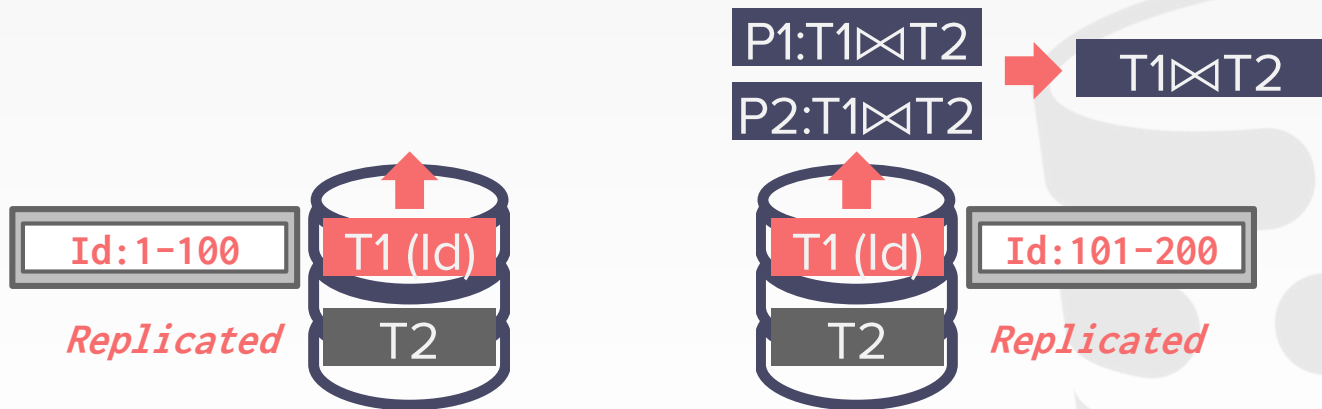
```
SELECT * FROM T1, T2
WHERE T1.id = T2.id
```

# SCENARIO #2

Tables are partitioned on the join attribute. Each node performs the join on local data and then sends to a node for coalescing.

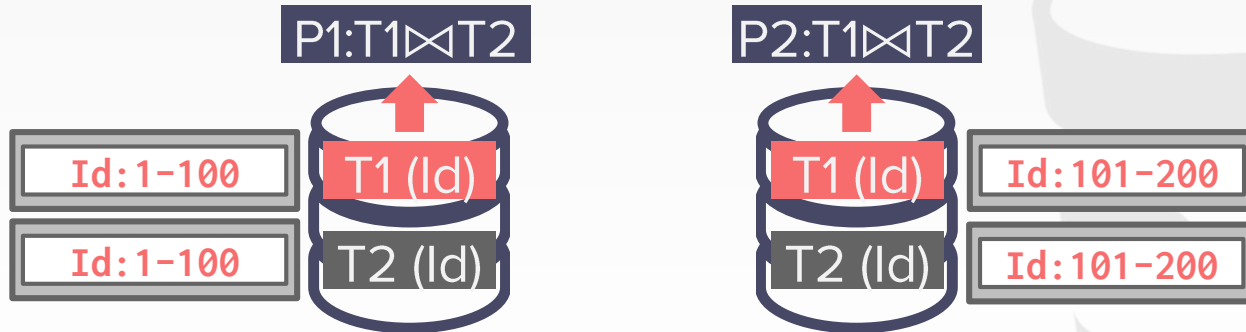```
SELECT * FROM T1, T2
  WHERE T1.id = T2.id
```

# SCENARIO #2

Tables are partitioned on the join attribute. Each node performs the join on local data and then sends to a node for coalescing.
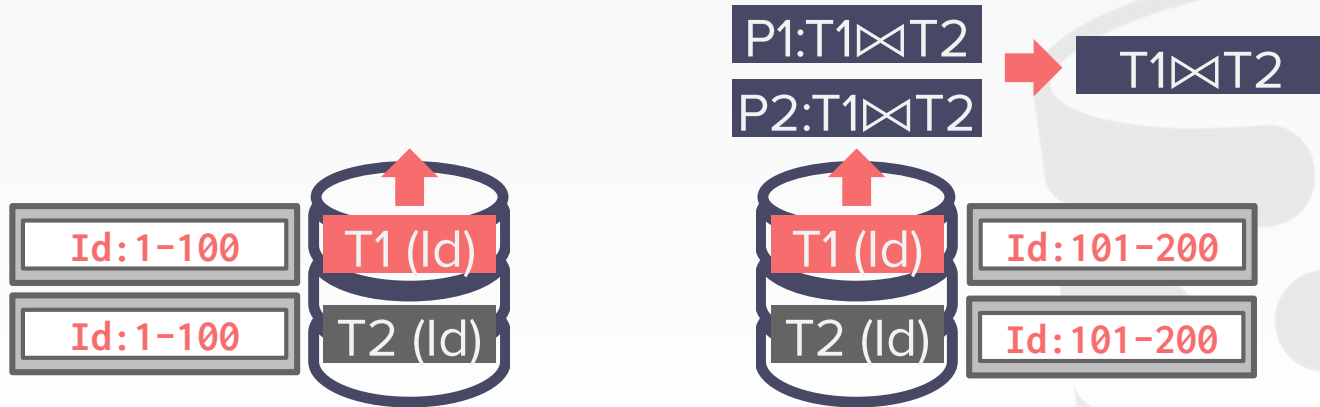
```
SELECT * FROM T1, T2
  WHERE T1.id = T2.id
```

# SCENARIO #2

Tables are partitioned on the join attribute. Each node performs the join on local data and then sends to a node for coalescing.

```
SELECT * FROM T1, T2
WHERE T1.id = T2.id
```

P1:T1⋈T2

P2:T1⋈T2

T1⋈T2

Id:1-100    T1 (Id)

Id:1-100    T2 (Id)

T1 (Id)    Id:101-200

T2 (Id)    Id:101-200

CARNEGIE MELLON
DATABASE GROUP

# SCENARIO #3

Both tables are partitioned on different keys. If one of the tables is small, then the DBMS **broadcasts** that table to all nodes.

```
SELECT * FROM T1, T2
 WHERE T1.id = T2.id
```



Id:1-100  |  T1 (Id)

Val:1-50  |  T2 (Val)

T1 (P2)  |  Id:101-200

T2 (Val)  |  Val:51-100

CARNEGIE MELLON
**DATABASE GROUP**

# SCENARIO #3

Both tables are partitioned on different keys. If one of the tables is small, then the DBMS **broadcasts** that table to all nodes.

```
SELECT * FROM T1, T2
  WHERE T1.id = T2.id
```

CARNEGIE MELLON
**DATABASE GROUP**

# SCENARIO #3

Both tables are partitioned on different keys. If one of the tables is small, then the DBMS **broadcasts** that table to all nodes.
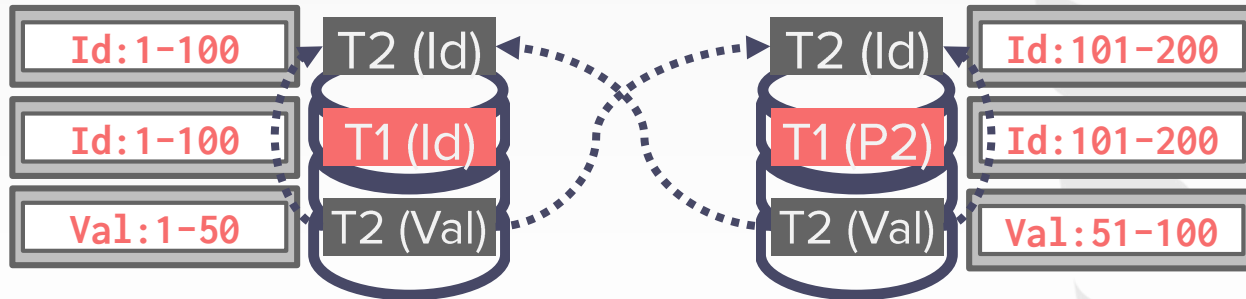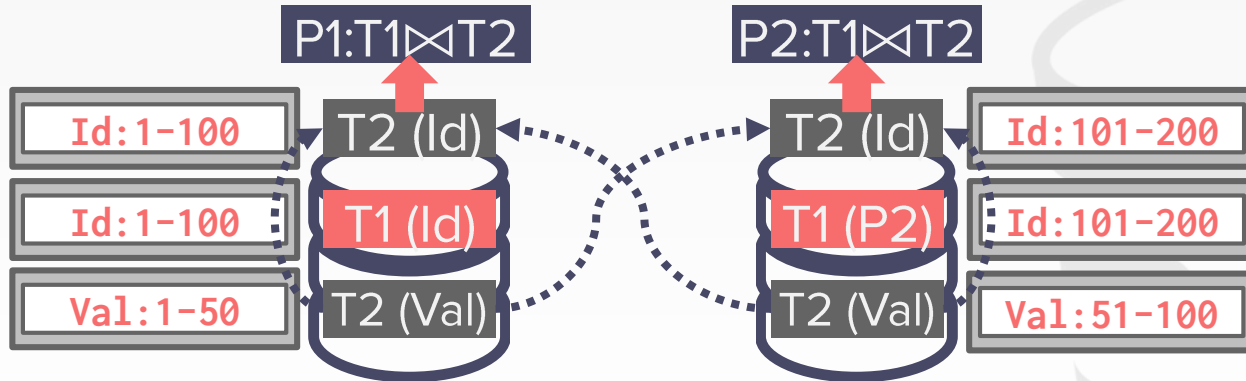
```
SELECT * FROM T1, T2
  WHERE T1.id = T2.id
```
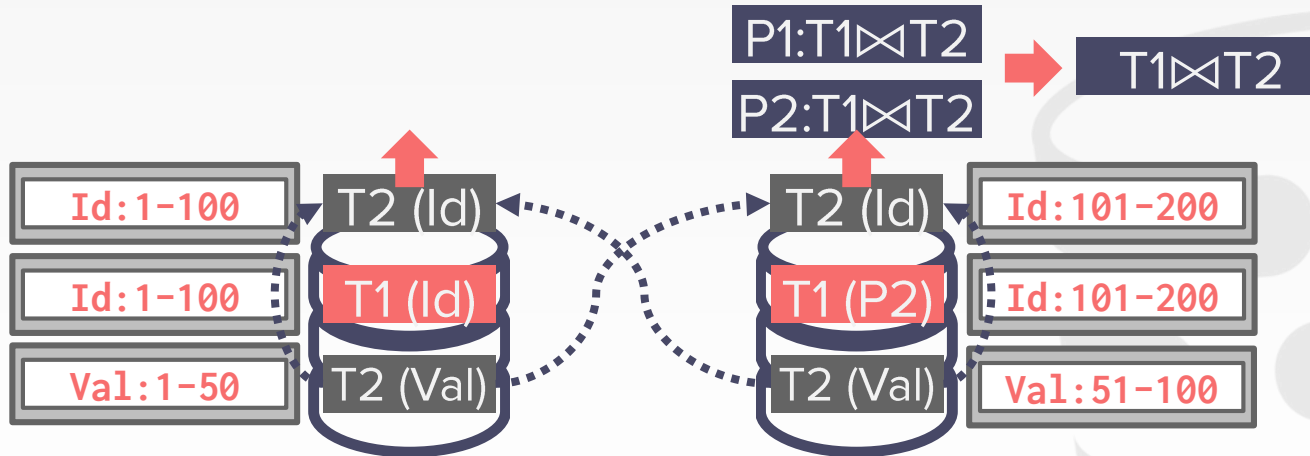
# SCENARIO #3

Both tables are partitioned on different keys. If one of the tables is small, then the DBMS **broadcasts** that table to all nodes.

```
SELECT * FROM T1, T2
 WHERE T1.id = T2.id
```

CARNEGIE MELLON
**DATABASE GROUP**

# SCENARIO #4

Both tables are not partitioned on the join key. The DBMS copies the tables by **reshuffling** them across nodes.

```
SELECT * FROM T1, T2
 WHERE T1.id = T2.id
```

| Name:A-M | T1 (Name) |
| Val:1-50 | T2 (Val) |

| T1 (Name) | Name:N-Z |
| T2 (Val) | Val:51-100 |

CARNEGIE MELLON
DATABASE GROUP

# SCENARIO #4

Both tables are not partitioned on the join key. The DBMS copies the tables by **reshuffling** them across nodes.

```
SELECT * FROM T1, T2
 WHERE T1.id = T2.id
```
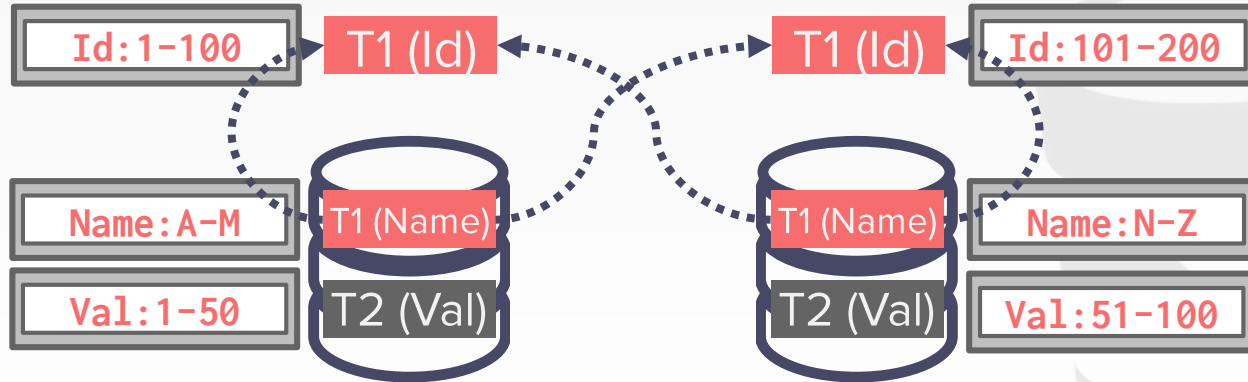
CARNEGIE MELLON
DATABASE GROUP

# SCENARIO #4

Both tables are not partitioned on the join key. The DBMS copies the tables by **reshuffling** them across nodes.

```
SELECT * FROM T1, T2
 WHERE T1.id = T2.id
```

CARNEGIE MELLON
DATABASE GROUP

# SCENARIO #4

Both tables are not partitioned on the join key. The DBMS copies the tables by **reshuffling** them across nodes.
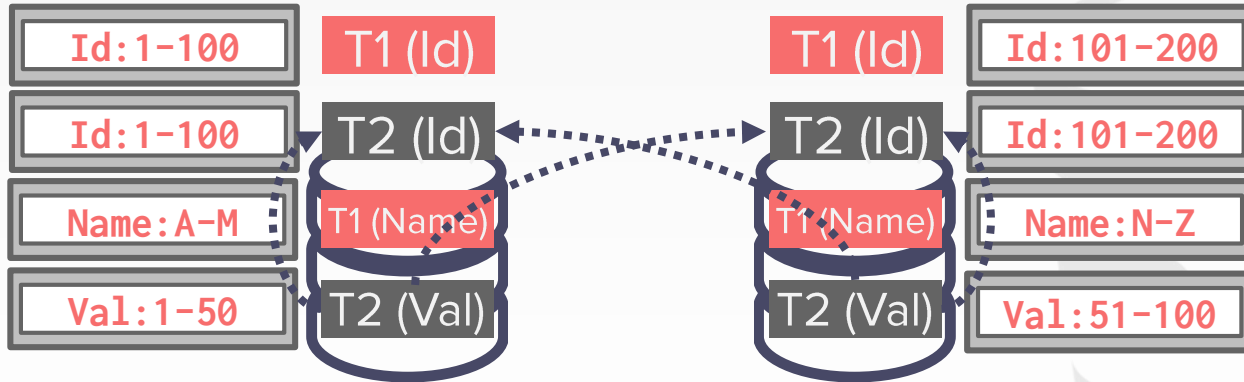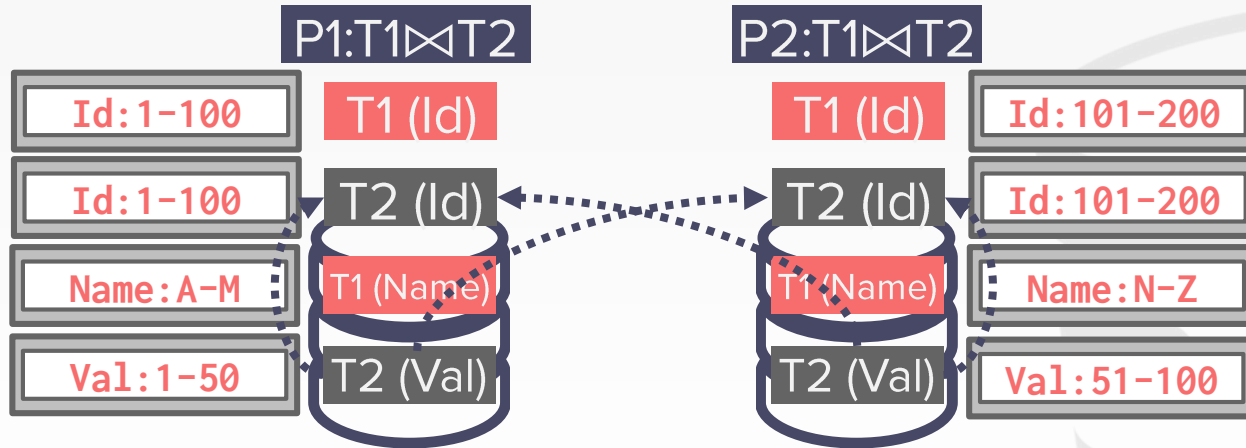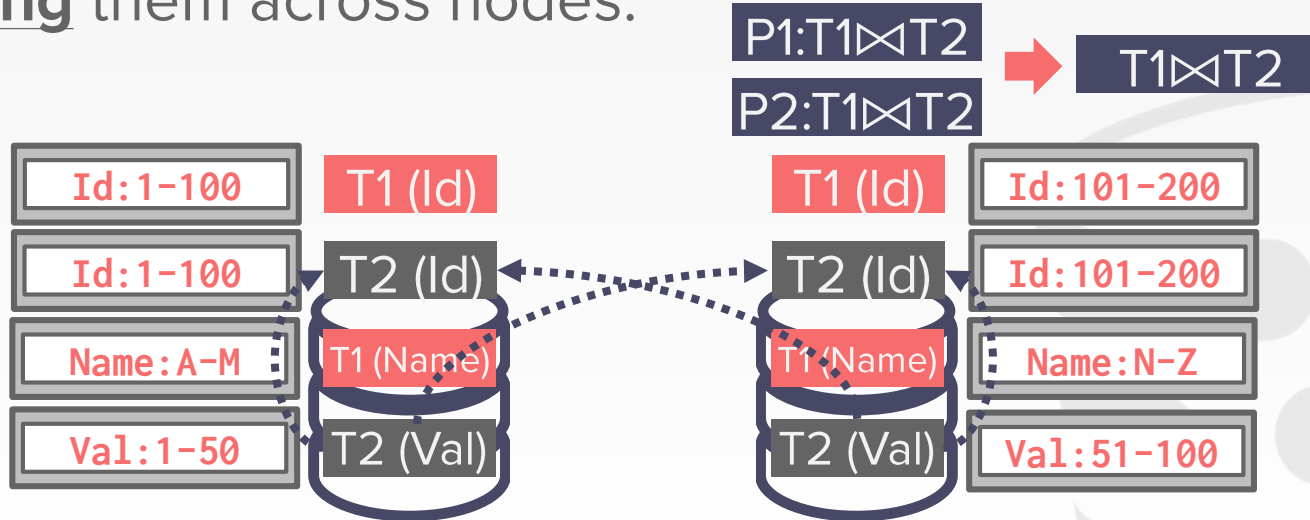
```
SELECT * FROM T1, T2
 WHERE T1.id = T2.id
```

# SCENARIO #4

Both tables are not partitioned on the join key. The DBMS copies the tables by **reshuffling** them across nodes.

```
SELECT * FROM T1, T2
  WHERE T1.id = T2.id
```

P1:T1⋈T2

P2:T1⋈T2 → T1⋈T2

| Id:1-100 | T1 (Id) | T1 (Id) | Id:101-200 |
| Id:1-100 | T2 (Id) | T2 (Id) | Id:101-200 |
| Name:A-M | T1 (Name) | T1 (Name) | Name:N-Z |
| Val:1-50 | T2 (Val) | T2 (Val) | Val:51-100 |

CARNEGIE MELLON
DATABASE GROUP

# CONCLUSION

Again, efficient distributed OLAP systems are difficult to implement.

Whenever possible, you want to push the query to the data rather than pull the data to the query.

CARNEGIE MELLON
DATABASE GROUP

# NEXT CLASS

NuoDB!

CARNEGIE MELLON
DATABASE GROUP