# Join Types - Big Data

March 8, 2025

# 1 Understanding Join Types in Big Data

## 1.1 Table of Contents

## 1.2 1. Introduction to Joins

In big data processing, a **join operation** combines rows from two or more tables based on a related column between them. Joins are fundamental operations in data processing and analysis, allowing us to connect information from different sources.

### 1.2.1 Why are Joins Important in Big Data?

- **Data Integration**: Combine data from multiple sources
- **Normalization Support**: Work with normalized data structures efficiently
- **Relationship Analysis**: Understand connections between different entities
- **Enrichment**: Add context and additional information to core datasets

### 1.2.2 Join Fundamentals

- **Join Key**: The column(s) used to match rows between tables
- **Join Condition**: The criteria that determines which rows are matched
- **Join Type**: Determines which rows are included in the result
- **Join Algorithm**: The method used to perform the join operation

```
 Table A                    Table B

key_column             key_column
column_a1                  column_b1
```

```
column_a2                      column_b2
```

```
              Joined Result

              key_column
              column_a1
              column_a2
              column_b1
              column_b2
```

## 1.3   2. Join Algorithms

Before diving into join types, it's important to understand the algorithms used to execute joins in big data systems. The choice of algorithm significantly impacts performance.

### 1.3.1   Shuffle Hash Join (Regular Join)

In a shuffle hash join (often called a regular join), data from both tables is shuffled across the network so that matching keys are processed by the same executor.

**Process**: 1. Both tables are shuffled and partitioned by the join key 2. Matching partitions are sent to the same executor 3. For each partition, a hash table is built for the smaller table 4. The larger table is scanned and matched against the hash table

```
  Table A                   Table B




        Shuffle by Join Key





  Partition              Partition
     A1                      B1




        Hash Join
```

```
        Results
```

**Pros**: - Works for any size of tables - Evenly distributes the processing

**Cons**: - High network shuffling cost - Expensive for large datasets

### 1.3.2 Broadcast Hash Join

In a broadcast hash join, the smaller table is broadcast (copied) to all executors processing the larger table, eliminating the need for shuffling.

**Process**: 1. The smaller table is broadcast to all executors 2. Each executor builds a hash table of the smaller table 3. Each partition of the larger table is processed locally 4. Rows are matched against the hash table

```
  Small                   Large
  Table A                 Table B




  Broadcast
  to all
  Executors




Copy A Copy A Copy A    Partition
                  B1



Join    Join   Join     Partition
A & B1 A & B2 A & B3       B2



                  Partition
                    B3
```

3

```
        Results
```

**Pros**: - Eliminates shuffling - Significantly faster for small-large table joins

**Cons**: - Works only when one table is small enough to fit in memory - Creates memory pressure if broadcast table is too large

### 1.3.3  Sort Merge Join

This algorithm first sorts both datasets by the join key and then merges them together.

**Process**: 1. Both tables are shuffled by the join key 2. Each partition is sorted by the join key 3. The sorted partitions are merged together

```
 Table A                Table B




 Shuffle &              Shuffle &
   Sort                   Sort




            Merge




            Results
```

**Pros**: - Works well for large datasets - More memory efficient than hash joins

**Cons**: - Sorting can be expensive - Not as fast as broadcast join for small-large table combinations

## 1.4  3. Types of Joins

Join types determine which records are included in the result set based on matching conditions. Let's explore the most common join types using example tables.

### 1.4.1  Example Tables

For our examples, we'll use two simple tables: `Customers` and `Orders`.

**Customers Table**

| customer_id | name | city |
|---|---|---|
| 101 | Alice | New York |
| 102 | Bob | Chicago |
| 103 | Charlie | Boston |
| 104 | David | Seattle |

**Orders Table**

| order_id | customer_id | amount | status |
|---|---|---|---|
| 1001 | 101 | 100.00 | Shipped |
| 1002 | 102 | 200.00 | Pending |
| 1003 | 101 | 150.00 | Delivered |
| 1004 | 105 | 300.00 | Shipped |

Note: Customer 103 and 104 have no orders, and order 1004 is for customer 105 who is not in the Customers table.

These tables will be used throughout the notebook to demonstrate different join types.

## 1.5  4. Inner Join

An **inner join** returns only the rows that have matching values in both tables. It's the most common type of join.

### 1.5.1  Visual Representation

```
    Table A              Table B




            Result
```

### 1.5.2  Example

Let's perform an inner join on the `Customers` and `Orders` tables:

```
SELECT c.customer_id, c.name, o.order_id, o.amount
FROM Customers c
INNER JOIN Orders o ON c.customer_id = o.customer_id
```

**Result**

| customer_id | name | order_id | amount |
|---|---|---|---|
| 101 | Alice | 1001 | 100.00 |
| 101 | Alice | 1003 | 150.00 |
| 102 | Bob | 1002 | 200.00 |

### 1.5.3 Characteristics

- **Inclusion**: Only rows with matching values in both tables
- **NULL Values**: Rows with NULL in the join columns are not included
- **Duplicate Handling**: If multiple matches exist, all combinations are returned
- **Use Cases**:
  - Finding relationships that exist in both tables
  - Combining related data when you only care about complete information

### 1.5.4 Real-world Scenario

Finding all customers who have placed orders, along with their order details. This excludes customers who haven't placed orders and orders from unregistered customers.

## 1.6 5. Outer Joins

Outer joins return all rows from one or both tables, regardless of whether there are matching values. There are three types: **LEFT**, **RIGHT**, and **FULL** outer joins.

### 1.6.1 Left Outer Join

A **left outer join** returns all rows from the left table, and the matched rows from the right table. When there's no match, NULL values are filled in for the right table's columns.

**Visual Representation**

```
    Table A            Table B
    (Left)             (Right)




            Result
```

**Example**

```sql
SELECT c.customer_id, c.name, o.order_id, o.amount
FROM Customers c
LEFT JOIN Orders o ON c.customer_id = o.customer_id
```

**Result**

| customer_id | name | order_id | amount |
|---|---|---|---|
| 101 | Alice | 1001 | 100.00 |
| 101 | Alice | 1003 | 150.00 |
| 102 | Bob | 1002 | 200.00 |
| 103 | Charlie | NULL | NULL |
| 104 | David | NULL | NULL |

**Use Cases**: - Finding all customers and their orders (if any) - Identifying customers who haven't placed orders

### 1.6.2 Right Outer Join

A **right outer join** returns all rows from the right table, and the matched rows from the left table. When there's no match, NULL values are filled in for the left table's columns.

**Visual Representation**

```
    Table A                 Table B
    (Left)                  (Right)




            Result
```

**Example**

```sql
SELECT c.customer_id, c.name, o.order_id, o.amount
FROM Customers c
RIGHT JOIN Orders o ON c.customer_id = o.customer_id
```

**Result**

| customer_id | name | order_id | amount |
|---|---|---|---|
| 101 | Alice | 1001 | 100.00 |
| 101 | Alice | 1003 | 150.00 |
| 102 | Bob | 1002 | 200.00 |
| NULL | NULL | 1004 | 300.00 |

**Use Cases**: - Finding all orders and their corresponding customers (if any) - Identifying orders placed by unregistered customers

### 1.6.3 Full Outer Join

A **full outer join** returns all rows from both tables. When there's no match, NULL values are filled in for the columns from the table without a match.

**Visual Representation**

Table A                 Table B

            Result

**Example**

```
SELECT c.customer_id, c.name, o.order_id, o.amount
FROM Customers c
FULL OUTER JOIN Orders o ON c.customer_id = o.customer_id
```

**Result**

| customer_id | name | order_id | amount |
|---|---|---|---|
| 101 | Alice | 1001 | 100.00 |
| 101 | Alice | 1003 | 150.00 |
| 102 | Bob | 1002 | 200.00 |
| 103 | Charlie | NULL | NULL |
| 104 | David | NULL | NULL |
| NULL | NULL | 1004 | 300.00 |

**Use Cases**: - Getting a complete view of all data from both tables - Identifying data mismatches

between systems - Finding data that exists in only one of the tables

## 1.7  6. Semi Joins and Anti Joins

These join types are used for **filtering data** rather than combining columns from multiple tables.

### 1.7.1  Left Semi Join

A **left semi join** returns rows from the left table where there is a match with the right table, but does not include columns from the right table.

**Visual Representation**

```
    Table A                Table B




            Table A Rows
            (with match)
```

**Example**

```sql
SELECT c.customer_id, c.name
FROM Customers c
LEFT SEMI JOIN Orders o ON c.customer_id = o.customer_id
```

**Result**

| customer_id | name |
|-------------|------|
| 101 | Alice |
| 102 | Bob |

**Use Cases**: - Finding customers who have placed at least one order - Filtering records based on existence in another table - Equivalent to `EXISTS` subquery in SQL

### 1.7.2  Left Anti Join

A **left anti join** returns rows from the left table where there is no match with the right table.

**Visual Representation**

```
    Table A              Table B



                \



            Table A Rows
          (without match)
```

**Example**

```sql
SELECT c.customer_id, c.name
FROM Customers c
LEFT ANTI JOIN Orders o ON c.customer_id = o.customer_id
```

**Result**

| customer_id | name |
| --- | --- |
| 103 | Charlie |
| 104 | David |

**Use Cases**: - Finding customers who have not placed any orders - Identifying missing relationships - Equivalent to `NOT EXISTS` subquery in SQL

## 1.8  7. Cross Join

A **cross join** returns the Cartesian product of the two tables, combining each row from the first table with each row from the second table.

**Visual Representation**

```
    Table A              Table B




           ×



         All possible
         combinations
```

**Example**

```sql
SELECT c.name, p.product_name
FROM Customers c
CROSS JOIN Products p
```

Assuming a `Products` table with 3 products, the result would combine each customer with each product, resulting in 12 rows (4 customers × 3 products).

**Characteristics**: - **Result Size**: Multiplies the number of rows (m × n) - **Join Condition**: No join condition is required - **Performance Impact**: Can be very resource-intensive for large tables

**Use Cases**: - Generating all possible combinations (e.g., product configurations) - Creating test data - Data permutation analysis

## 1.9 8. Performance Considerations

Join performance is a critical concern in big data processing. Here are key factors that affect join performance:

### 1.9.1 Impact of Data Size

| Join Type | Small Tables | Large Tables |
|-----------|--------------|--------------|
| Inner Join | Fast | Slower due to shuffling |
| Outer Join | Fast | Slower due to shuffling |
| Cross Join | Fast | Very slow (Cartesian product) |

### 1.9.2 Join Algorithm Selection

| Algorithm | Best For | Worst For |
|-----------|----------|-----------|
| Shuffle Hash Join | Large tables | Small-large table joins |
| Broadcast Hash Join | Small-large table joins | Large tables |
| Sort Merge Join | Large tables | Small tables |

### 1.9.3 Optimization Tips

- **Broadcast Small Tables**: Use broadcast joins for small-large table combinations
- **Partitioning**: Partition tables on join keys to reduce shuffling
- **Bucketing**: Use bucketing to optimize joins on large tables
- **Filter Early**: Apply filters before joins to reduce data size
- **Compression**: Use compressed file formats like Parquet or ORC

## 1.10 9. Best Practices

To get the most out of joins, follow these best practices:

### 1.10.1 Choosing Join Types

- **Inner Join**: When you only need matching rows

- **Left Join**: When you need all rows from the left table
- **Right Join**: When you need all rows from the right table
- **Full Join**: When you need all rows from both tables
- **Semi Join**: When you need to filter rows based on existence
- **Anti Join**: When you need to filter rows based on absence
- **Cross Join**: When you need all possible combinations

### 1.10.2 Join Key Selection

- **Primary Key**: Use primary keys for joins when possible
- **Indexed Columns**: Use indexed columns for faster joins
- **Low Cardinality**: Avoid high-cardinality columns for joins

### 1.10.3 Performance Optimization

- **Broadcast Small Tables**: Use broadcast joins for small-large table combinations
- **Partitioning**: Partition tables on join keys to reduce shuffling
- **Bucketing**: Use bucketing to optimize joins on large tables
- **Filter Early**: Apply filters before joins to reduce data size
- **Compression**: Use compressed file formats like Parquet or ORC

## 1.11 10. Summary

Joins are a fundamental operation in big data processing, allowing us to combine data from multiple sources. Here's a recap of the key points:

### 1.11.1 Key Join Types

| Join Type | Description | Use Case |
| --- | --- | --- |
| Inner Join | Returns matching rows | Finding relationships |
| Left Join | Returns all rows from left table | Including unmatched rows |
| Right Join | Returns all rows from right table | Including unmatched rows |
| Full Join | Returns all rows from both tables | Complete data view |
| Semi Join | Filters rows based on existence | Filtering data |
| Anti Join | Filters rows based on absence | Filtering data |
| Cross Join | Returns Cartesian product | Generating combinations |

### 1.11.2 Best Practices

- Choose the right join type for your use case
- Optimize join performance with partitioning and bucketing
- Use broadcast joins for small-large table combinations
- Filter data early to reduce join size
- Use compressed file formats like Parquet or ORC

By following these principles and best practices, you can significantly improve the performance and efficiency of your big data processing tasks.