**Spark Join Strategies Analysis**

The main strategies are:

**1. Shuffle Hash Join**

**Description**: Both datasets are shuffled by the join key, then partitions with the same key are hashed and joined.

**Characteristics**:

- Requires shuffling of data across the network
- Builds a hash table on the smaller dataset for each partition
- Effective when joining a large dataset with a smaller one
- Performance can degrade with larger datasets due to shuffle overhead

**Use cases**:

- When one of the tables is small enough to fit in memory (controlled by spark.sql.autoBroadcastJoinThreshold)
- When one of the tables is small enough to fit in memory but too large to broadcast
- When sort-merge join is not possible (e.g., with non-equality joins)

**2. Broadcast Hash Join**

**Description**: Spark broadcasts the smaller dataset to all executors, allowing each partition of the larger dataset to join efficiently.

**Characteristics**:

- Eliminates the need for shuffling
- Extremely efficient when one dataset is small enough to fit in memory
- Default threshold for broadcasting is 10MB (configurable via spark.sql.autoBroadcastJoinThreshold)
- Broadcasting fails when the dataset exceeds available executor memory

**Use cases**:

- When one relation is small enough to fit in the executor memory
- When the join condition uses equality operators

**3. Sort-Merge Join**

**Description**: Both datasets are sorted by the join key, then merged together.

**Characteristics**:

- Works well for large datasets that don't fit in memory
- Requires sorting both datasets, which can be expensive

- Efficient for large-to-large table joins
- Default join strategy in Spark for equi-joins since Spark 2.3

**Use cases**:

- When both datasets are large
- When the join uses equality comparison
- When spark.sql.join.preferSortMergeJoin is true (default in modern Spark)

## 4. Cartesian Join

**Description**: Every row from the first dataset is joined with every row from the second dataset.

**Characteristics**:

- Extremely expensive in terms of computation
- Results in m × n rows where m and n are the sizes of the input datasets
- Rarely used in production environments

**Use cases**:

- When no join condition is specified unless explicitly triggered using .crossJoin()
- When the join condition doesn't allow more efficient strategies

## 5. Broadcast Nested Loop Join

**Description**: Used for non-equi joins, where one dataset (usually smaller) is broadcast to all executors.

**Characteristics**:

- More efficient than a regular cartesian product
- Used for complex join conditions that aren't equality-based
- Considerably more expensive than hash or sort-merge joins

**Use cases**:

- When Spark cannot use hash or sort-merge joins
- When the join condition isn't an equality
- When one dataset is small enough to broadcast

**Strategy for joining Parquet File 1 and 2 in Spark Dataframe**

When joining Dataset A (~1 million rows) with Dataset B (10,000 rows) on geographical_location_oid, here's why **Broadcast Hash Join** is the best fit:

1. **Dataset Sizes**:

   o   Dataset A is large (~1 million rows), while Dataset B is much smaller (10,000 rows). The small size of Dataset B makes it a perfect candidate for broadcasting as it can easily fit in memory on each executor

2. **Performance:**

   o   By broadcasting Dataset B, we eliminate the need for shuffling, which makes the join faster and reduces network overhead.

3. **Ease of Implementation:**

   o   Broadcasting is simple to implement and saves computational resources. Plus, it's very efficient when one dataset fits in memory.

4. **Data Skew**:

   o   Even if there's some skew in Dataset A, the fact that Dataset B is broadcasted makes this less of a concern because it doesn't depend on how the data is distributed across nodes.

**Conclusion**

For joining Parquet File 1 and Parquet File 2, Broadcast Hash Join is the most efficient strategy. It optimizes memory and network usage, speeds up the join, and handles data skew effectively. With the small size of Dataset B, this approach will ensure the best performance and ease of implementation.

The broadcast hash join will follow a simple implementation that looks like this:

```scala
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.functions._
import org.apache.spark.sql.expressions.Window
import org.apache.spark.sql.DataFrame

// Read parquet files
val dataA = spark.read.parquet("path/to/parquet1")
val dataB = spark.read.parquet("path/to/parquet2")

// Deduplicate detection events and select only needed columns
val dedupedDF = dataA
  .dropDuplicates("detection_oid")
  .select("geographical_location_oid", "item_name")

// Get item counts by location
val itemCountDF = dedupedDF
  .groupBy("geographical_location_oid", "item_name")
  .count()

// Create window specification for ranking items within each location
val windowSpec = Window
  .partitionBy("geographical_location_oid")
  .orderBy(desc("count"))

// Add ranking information
val rankedDF = itemCountDF
  .withColumn("item_rank", row_number().over(windowSpec).cast("string"))
  .filter(col("item_rank") <= lit(5)) // Select top 5 items per location

// Join with location data to get location names
val combinedDF = rankedDF.join(
  broadcast(dataB),
  rankedDF("geographical_location_oid") === dataB("geographical_location_oid")
)

// Select final columns for output
val resultDF = combinedDF
  .select(
    col("geographical_location"),
    col("item_rank"),
    col("item_name")
  )
```

**Handling data skew**

For handling data skew we can perform post-join aggregations where skew might still affect performance.

```scala
import org.apache.spark.sql.functions._
import org.apache.spark.sql.expressions.Window

// First identify skewed locations by counting records per location
val locationCounts = dataA
  .groupBy("geographical_location_oid")
  .count()
  .cache()

// Find the most skewed location
val skewedLocation = locationCounts
  .orderBy(desc("count"))
  .first()
  .getAs[Long]("geographical_location_oid")

// Define the number of partitions to use for salting
val numPartitions = 10

// Split data into skewed and non-skewed paths
val skewedData = dataA.filter(col("geographical_location_oid") === skewedLocation)
val normalData = dataA.filter(col("geographical_location_oid") =!= skewedLocation)

// Process normal data using the standard approach
val normalResults = normalData
  .dropDuplicates("detection_oid")
  .groupBy("geographical_location_oid", "item_name")
  .count()

// Process skewed data using salting technique
val saltedResults = skewedData
  // Add deterministic salt based on detection_oid
  .withColumn("salt", mod(col("detection_oid"), lit(numPartitions)))
  // Deduplicate by detection_oid
  .dropDuplicates("detection_oid")
  // First level aggregation with salt to distribute the work
  .groupBy("geographical_location_oid", "item_name", "salt")
  .count()
  // Second level aggregation to combine salted results
  .groupBy("geographical_location_oid", "item_name")
  .agg(sum("count").as("count"))

// Combine results from both paths
val combinedResults = normalResults.union(saltedResults)

// Create window spec for ranking
val windowSpec = Window
  .partitionBy("geographical_location_oid")
  .orderBy(desc("count"))

// Generate final ranked output
val rankedItemsDF = combinedResults
  .withColumn("item_rank", row_number().over(windowSpec).cast("string"))
  .filter(col("item_rank").leq(lit(5))) // Top 5 items per location
  .join(dataB, "geographical_location_oid")
  .select(
    col("geographical_location"),
    col("item_rank"),        You, 1 minute ago • Uncommitted changes
    col("item_name")
  )
```

**Conclusion**

The Broadcast Hash Join strategy would be the most efficient approach for this specific dataset. This strategy optimizes both computational and network resources while effectively handling any data skew present in the geographical locations.