# Sorting Strategies in Spark

**Sort-Merge Join**

Sort Merge is the default join strategy in Spark.  It sorts both DataFrames on the join keys and then merges them together.   Sort-Merge Join is particularly efficient for large datasets in Spark. It works well when both datasets are too large to fit in memory and need to be sorted and merged. Here are some key points about its efficiency:

- **Scalability**: Sort-Merge Join can handle datasets that are several terabytes in size.  It leverages Spark's ability to distribute and parallelize the sorting and merging process across the cluster.
- **Partitioning**: The efficiency of Sort-Merge Join depends on how well the data is partitioned. If the data is evenly distributed across partitions, the join operation will be more efficient
- **Sorting**: Both datasets are sorted on the join keys before the merge. This sorting step can be resource-intensive, but it ensures that the merge operation is efficient
- **Shuffling**: Sort-Merge Join involves shuffling data across the cluster to ensure that rows with the same join keys end up in the same partition. While shuffling can be expensive, Spark optimizes this process to minimize overhead

**Broadcast Join**

Broadcast Join is a powerful technique in Spark for optimizing join operations when one dataset is significantly smaller than the other. Here are some key points about its efficiency:

- **Size Threshold**: Spark uses the spark.sql.autoBroadcastJoinThreshold property to determine whether to use a broadcast join. By default, this threshold is set to 10MB.  You can adjust this threshold based on your cluster's memory capacity.
- **Memory Usage**: Broadcast joins are efficient when the smaller dataset can fit into the memory of each executor. This minimizes shuffling and network overhead, making the join operation faster.
- **Performance**: Broadcasting the smaller dataset to all nodes in the cluster allows each node to perform the join locally, reducing the need for data movement across the network.  This can significantly speed up the join operation, especially when dealing with large datasets.

**Shuffle Hash Join**

Shuffle Hash Join is another strategy in Spark for joining large datasets. It can be efficient under certain conditions, particularly when the data is already partitioned by the join keys. Here are some key points about its efficiency:

- **Partitioning**: Shuffle Hash Join partitions both DataFrames by the join keys and then performs a hash join on the partitions. This can be efficient if the data is already partitioned by the join keys, reducing the need for additional shuffling.
- **Data Distribution**: It works well when the data is evenly distributed across partitions. If there is data skew, it can lead to some partitions being much larger than others, which can slow down the join operation.
- **Resource Usage**: Shuffle Hash Join can be resource-intensive because it involves shuffling data across the cluster. However, it can be optimized by ensuring that the data is well-partitioned and by using appropriate configurations.

**Recommended Strategy**

For joining Parquet File 1 and 2, I recommend using the **Broadcast Join** strategy as File 2 (DataSet B - Geographical Location Master) is small enough to fit in memory. This will minimize shuffling and speed up the join operation.