Delving deeper into Apache Impala, an essential tool for real-time querying in Hadoop environments, let's explore its advanced usage, including complex steps, code, and scripts.

**Deep Dive into Apache Impala**

**Architecture Overview**

- **Components**: Impala consists of the Impala Daemon (impalad), the StateStore (statestored), and the Catalog Service (catalogd).
  - **Impala Daemon**: Runs on each Hadoop node and is responsible for reading and writing data to/from HDFS or HBase.
  - **StateStore**: Monitors the health of Impala Daemons and coordinates actions across the cluster.
  - **Catalog Service**: Manages metadata and broadcasts changes to all Impala Daemons.

**Complex Impala Query**

Impala supports advanced SQL features. Here's an example of a complex query involving joins, aggregations, and window functions:

```sql
SELECT a.user_id, a.transaction_id, a.amount,
       SUM(a.amount) OVER (PARTITION BY a.user_id ORDER BY a.transaction_date
                              ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW) AS running_tota
FROM transactions a
JOIN users b ON a.user_id = b.id
WHERE b.region = 'North America'
AND a.transaction_date BETWEEN '2024-01-01' AND '2024-12-31';
```

This query calculates the running total of transaction amounts for each user in North America for the year 2024.

**Impala Table Design and Partitioning**

Creating an efficient table in Impala involves partitioning and choosing the right file format:

```sql
CREATE TABLE large_sales (
    sale_id INT,
    product_id INT,
    quantity_sold INT,
    sale_date DATE,
    region STRING
)
PARTITIONED BY (year INT, month INT)
STORED AS PARQUET;
```

This table is partitioned by year and month, which helps in efficient querying, especially for large datasets.

### Impala Performance Tuning

- **Partition Pruning**: Effective use of partitioning in table design can significantly speed up queries.
- **Join Strategies**: Impala offers several join algorithms. Use the right join strategy (broadcast join, partitioned join) based on your data size and distribution.
- **Memory Management**: Configure memory usage per query to optimize performance, especially for large datasets.

### Impala and HBase

Impala can query HBase tables. Here's how you can map an HBase table in Impala:

```
CREATE EXTERNAL TABLE hbase_sales (
    key STRING,
    sale_date STRING,
    amount DOUBLE
)
STORED BY 'org.apache.hadoop.hive.hbase.HBaseStorageHandler'
WITH SERDEPROPERTIES ("hbase.columns.mapping" = ":key,cf:sale_date,cf:amount")
TBLPROPERTIES("hbase.table.name" = "sales");
```

This allows Impala to directly query data stored in HBase.

### Example: Complex Data Analysis Script

Suppose you need to analyze sales data. A complex script in Impala could involve multiple steps:

1. **Data Aggregation**:

   ```
   CREATE TABLE sales_summary AS
   SELECT region, SUM(amount) as total_sales
   FROM sales_data
   GROUP BY region;
   ```

2. **Data Join and Analysis**:

   ```
   SELECT a.region, a.total_sales, b.avg_price
   FROM sales_summary a
   JOIN (SELECT region, AVG(price) as avg_price FROM product_data GROUP BY region) b
   ON a.region = b.region;
   ```

3. **Temporary Views for Complex Analysis**:

```sql
CREATE VIEW temp_view AS
SELECT user_id, COUNT(*) as total_transactions
FROM transactions
WHERE transaction_date > '2024-01-01'
GROUP BY user_id;

SELECT * FROM temp_view WHERE total_transactions > 5;
```

**Note**

- Ensure your Impala cluster is properly sized and configured for the workload.
- Regularly update table statistics for optimal query performance.
- Use the `EXPLAIN` statement to understand and optimize query plans.