**Name: - Muhammad Huzaifa Waseem (2303-KHI-DEG-021)**
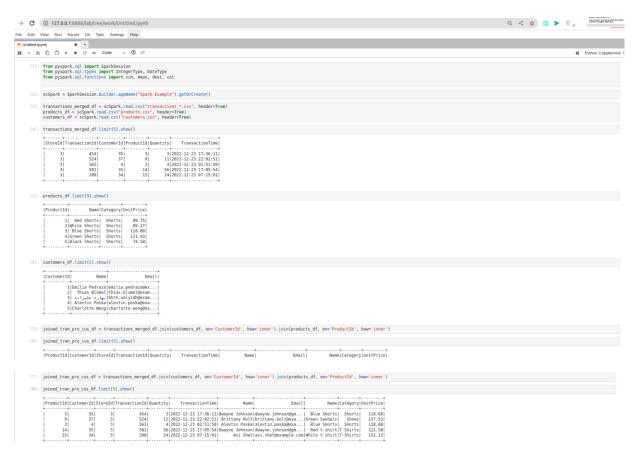
**Pair Partner 1: - Muhammad Faizan Rafique (2303.005.KHI.DEG)**

**Pair Partner 2: - Syed Muhammad Hammad Irshad(2303.KHI.DEG.032)**

**UNIT 5.1:**

# Assignment

Based on the data contained in *tasks/4_data_pipelines/day_1_introduction/daily_assignment/data* directory, use PySpark to read, filter and join the data from CSV files and answer the following questions:



- What are the daily total sales for the store with id 1?

```
[9]: # Convert TransactionTime column to DateType
     df = joined_tran_pro_cus_df.withColumn("TransactionTime", joined_tran_pro_cus_df["TransactionTime"].cast(DateType()))
```

```
[10]: # Convert Quantity and UnitPrice columns to appropriate data types
      df = df.withColumn("Quantity", df["Quantity"].cast(IntegerType()))
      df = df.withColumn("UnitPrice", df["UnitPrice"].cast(IntegerType()))

      # Calculate total sales for each transaction
      sales_df = df.withColumn("Sales", df["Quantity"] * df["UnitPrice"])

      # Group by TransactionTime and calculate daily total sales
      daily_total_sales_StoreId_1 = sales_df.filter(df["StoreId"] == 1).groupBy("TransactionTime").agg(sum("Sales").alias("TotalSales"))

      # Show the result
      daily_total_sales_StoreId_1.show()

      +---------------+----------+
      |TransactionTime|TotalSales|
      +---------------+----------+
      |     2022-12-23|     41070|
      +---------------+----------+
```

- ## What are the mean sales for the store with id 2?

```
[11]: # Group by TransactionTime and calculate daily total sales
      mean_sales_StoreId_2 = sales_df.filter(df["StoreId"] == 2).agg(mean("Sales").alias("TotalSales"))

      # Show the result
      mean_sales_StoreId_2.show()

      +----------------+
      |      TotalSales|
      +----------------+
      |511.921568627451|
      +----------------+
```

- ## What is the email of the client who spent the most when summing up purchases from all of the stores?

```
[12]: # Group by CustomerId and calculate the total purchase amount for each customer
      customer_total_purchase = sales_df.groupBy("CustomerId").agg(sum("Sales").alias("TotalPurchase"))
```

```
[13]: # Sort the data in descending order of the total purchase amount
      customer_total_purchase = customer_total_purchase.orderBy(desc("TotalPurchase")).limit(1)
```

```
[14]: # Get the email of the customer who spent the most
      most_spent_customer_email =customer_total_purchase.join(customers_df, on='CustomerId', how='inner').select("CustomerId", "TotalPurchase", "Email")
      most_spent_customer_email.show()

      +----------+-------------+-------------------+
      |CustomerId|TotalPurchase|              Email|
      +----------+-------------+-------------------+
      |        35|        10598|dwayne.johnson@gm...|
      +----------+-------------+-------------------+
```

- ## Which 5 products are most frequently bought across all stores?

```
[15]: frequently_bought_products = joined_tran_pro_cus_df.groupBy("ProductId", products_df["Name"]).agg(sum("Quantity").alias("TotalQuantity")).orderBy(desc("TotalQuantity")).limit(5)
      frequently_bought_products.show()

      +---------+-------------+-------------+
      |ProductId|         Name|TotalQuantity|
      +---------+-------------+-------------+
      |       14|  Red t-shirt|         82.0|
      |       24|   Blue Jeans|         77.0|
      |       15|White t-shirt|         76.0|
      |        5| Black Shorts|         75.0|
      |       19| Green jacket|         74.0|
      +---------+-------------+-------------+
```