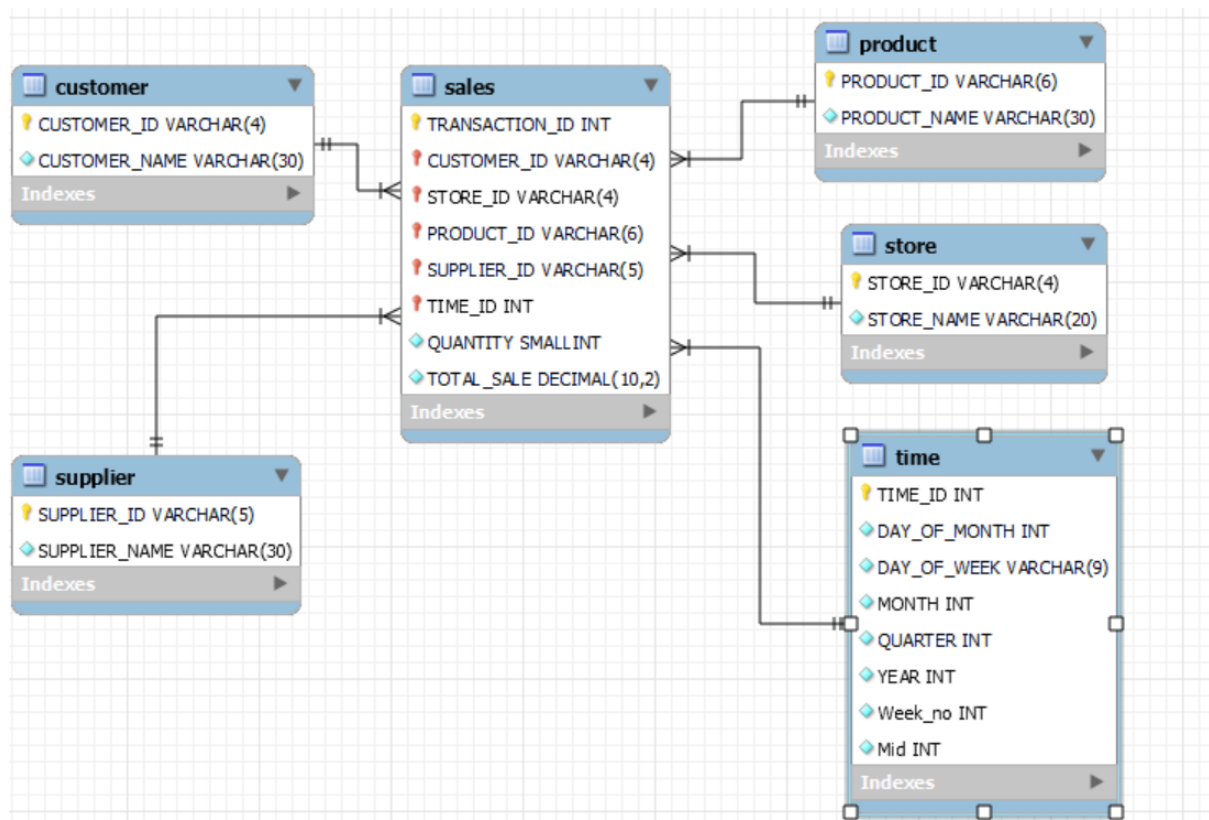# Project Overview

This project is to design and create such a data warehouse system Metro Pakistan | Biggest Wholesale & Supermarket Pakistan. The store has thousands of customers and therefore it is important for the store to online analyze the shopping behavior of their customers. Based on that the store can optimize their selling techniques e.g. giving promotions on different products. Due to the data from one-year transaction data.Transaction table 10000 records and master table 100 products. Program write for Extract Transform Load (ETL) in JAVA. When joining the transaction data and master data in the Meshjoin algorithm. The Start schema is used to model the data structure in the data warehouse first design the star schema. The start schema has five dimensions customer, store, Time, Product, Supplier, and Fact table Sale have two measures Total Sold and Total_Sale on the basis of the transaction. After Extracting and Transforming the data then loaded the data in the Data warehouse. Perform the Online Analytical Processing (OLAP) query.

# Schema for DW

The Start schema is used to model the data structure in the data warehouse. Star schema has five dimensions customer, Store, Product, Time, and Supplier. Fact sale table has two-measure Total sold and Total Sale.

# MESHJOIN Algorithm

The program dynamically loads the partition size and the stream partition's records size. The partition size of Meshjoin-partition is 10 and the stream partition size is 50. Add transaction records in an input queue. Enqueue the partition when it gets full. Add all new records into the record list. Load the next master partition data. Compare every record product id with master partition product id. Put a new record into the output queue. Last dequeue the last.

# Shortcomings in Mesh Join

1. Meshjoin depends on the partitions for the internal queue for The stream data and the number of changes needed for the disk based relation In memory . This dependency hampers the optimal distribution of memory between the Join components. Normally the size of the disk-buffer differs with the size of the disk-Based relation and its not necessary.

2. Stream-based algorithm Mesh Join (meshjoin) has been proposed to amortize disk Access over the fast stream. Meshjoin makes no assumptions about the data distribution. In Actual world applications, however, skewed distributions can be found, e.g, some products are sold more often than the remaining products. The question that comes to mind is, how much does the performance of mesh join loose by not adapting to data Skew.

3 The main problem in stream-relation joins is the different nature of inputs; stream data is Fast and busty, where the disk-based relation is not as fast because of the high disk I/O cost.

# Lessons Learned

- I learned java language programming and Logic builds Extraction Transform and Load (ETL) import data into Data warehouse.
- A real-time ETL program is to be implemented in reality.
- I learned OLAP querys Roll up , drill down

# DW analysis

## Q1 Present total sales of all products supplied by each supplier with respect to quarter and month

```
1 •   select l.SUPPLIER_NAME,t.QUARTER,t.MONTH,sum(s.TOTAL_SALE)
2     from data_warehouse_metro.sales s join data_warehouse_metro.supplier l
3     on(l.SUPPLIER_ID=s.SUPPLIER_ID)
4     join data_warehouse_metro.time t
5     on(t.TIME_ID=s.TIME_ID)
6     GROUP BY l.SUPPLIER_NAME,t.QUARTER,t.MONTH WITH ROLLUP
7     order by l.SUPPLIER_NAME,t.QUARTER,t.MONTH;
8
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: A

| SUPPLIER_NAME | QUARTER | MONTH | sum(s.TOTAL_SALE) |
|---|---|---|---|
| NULL | NULL | NULL | 717674.19 |
| 3Com Corp | NULL | NULL | 40192.94 |
| 3Com Corp | 1 | NULL | 9513.93 |
| 3Com Corp | 1 | 1 | 2987.05 |
| 3Com Corp | 1 | 2 | 2732.58 |
| 3Com Corp | 1 | 3 | 3794.30 |
| 3Com Corp | 2 | NULL | 9868.42 |

## Q2 Present total sales of each product sold by each store. The output should be organised store wise and then product wise under each store.

SQL File 8*    SQL File 9* ×

Limit to 50000 rows

```
11 •   select Store.STORE_NAME,Product.PRODUCT_NAME,sum(Sale.TOTAL_SALE)
12     from data_warehouse_metro.sales Sale join data_warehouse_metro.store Store
13     on(Store.STORE_ID=Sale.STORE_ID)
14     join data_warehouse_metro.product Product
15     on(Product.PRODUCT_ID=Sale.PRODUCT_ID)
16     GROUP BY Store.STORE_NAME,Product.PRODUCT_NAME  WITH ROLLUP
17     order by Store.STORE_NAME,Product.PRODUCT_NAME;
18
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: A

| STORE_NAME | PRODUCT_NAME | sum(Sale.TOTAL_SALE) |
|---|---|---|
| Albany | Apples | 581.44 |
| Albany | Applesauce | 1688.10 |
| Albany | Asparagus | 527.25 |
| Albany | Avocados | 717.66 |
| Albany | Bagels | 402.93 |
| Albany | Baked beans | 444.54 |
| Albany | Bananas | 681.05 |

Result 4    Result 5 ×

Result Grid

Form Editor

Read Only

## Q3 Find the 5 most popular products sold over the weekends.

```
20 ● ⊖ WITH sold AS (select p.PRODUCT_NAME,t.Week_no ,sum(s.QUANTITY),ROW_NUMBER() OVER (PARTITION BY t.Week_no
21     ⌐ ORDER BY sum(s.QUANTITY) DESC) top
22       from data_warehouse_metro.sales s join data_warehouse_metro.time t
23       on(t.TIME_ID=s.TIME_ID)
24       join data_warehouse_metro.product p
25       on(p.PRODUCT_ID=s.PRODUCT_ID)
26       GROUP BY p.PRODUCT_NAME,t.Week_no
27     ⌐ order by t.Week_no ,sum(s.QUANTITY) desc)
28       SELECT * FROM sold
29       WHERE top <= 5;
```

**Result Grid** | Filter Rows: [          ] | Export: | Wrap Cell Content: ‡A

| PRODUCT_NAME | Week_no | sum(s.QUANTITY) | top |
|---|---|---|---|
| Pickles | 2 | 31 | 3 |
| Tea | 2 | 30 | 4 |
| Bananas | 2 | 29 | 5 |
| Garlic | 3 | 28 | 1 |
| Oregano | 3 | 27 | 2 |

Result 6 ×                                                    ❶ Read Only   Cont

## Q4 Present the quarterly sales of each product for year 2016 using drill down query concept. Note: each quarter sale must be a column.

```
 1 ●    select p.PRODUCT_NAME,
 2        CASE t.QUARTER WHEN 1 THEN sum(s.TOTAL_SALE) ELSE NULL END as 'Quater_1',
 3     ⊖  CASE t.QUARTER WHEN 2 THEN sum(s.TOTAL_SALE)
 4        ELSE NULL
 5     ⌐  END as 'Quater_2',
 6     ⊖  CASE t.QUARTER WHEN 3 THEN sum(s.TOTAL_SALE)
 7        ELSE NULL
 8     ⌐  END as 'Quater_3',
 9     ⊖  CASE t.QUARTER WHEN 4 THEN sum(s.TOTAL_SALE)
10        ELSE NULL
11     ⌐  END as 'Quater_4'
12        from data_warehouse_metro.sales s join data_warehouse_metro.time t
13        on(t.TIME_ID=s.TIME_ID)
14        join data_warehouse_metro.product p
15        on(p.PRODUCT_ID=s.PRODUCT_ID)
16        GROUP BY p.PRODUCT_NAME,t.QUARTER
17        order by p.PRODUCT_NAME,t.QUARTER,sum(s.TOTAL_SALE);
18        |
```

Output

Action Output                            ▼

| # | Time | Action | Message |
|---|---|---|---|
| ✓ | 220 23:06:19 | WITH sold AS (select p.PRODUCT_NAME,t.Week_no ,sum(s.QUANTITY),ROW_... | 265 row(s) returned |
| ✓ | 221 23:11:04 | select p.PRODUCT_NAME, CASE t.QUARTER WHEN 1 THEN sum(s.TOTAL_SA... | 396 row(s) returned |

```
46        join data_warehouse_metro.product p
47        on(p.PRODUCT_ID=s.PRODUCT_ID)
48        GROUP BY p.PRODUCT_NAME,t.QUARTER
49        order by p.PRODUCT_NAME,t.QUARTER,sum(s.TOTAL_SALE);
50
51
52
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| PRODUCT_NAME | Quater_1 | Quater_2 | Quater_3 | Quater_4 |
|---|---|---|---|---|
| Apples | 1177.60 | NULL | NULL | NULL |
| Apples | NULL | 1096.64 | NULL | NULL |
| Apples | NULL | NULL | 920.00 | NULL |
| Apples | NULL | NULL | NULL | 1354.24 |
| Applesauce | 2101.85 | NULL | NULL | NULL |
| Applesauce | NULL | 2482.50 | NULL | NULL |
| Applesauce | NULL | NULL | 2813.50 | NULL |
| Applesauce | NULL | NULL | NULL | 2465.95 |
| Asparagus | 612.75 | NULL | NULL | NULL |
| Asparagus | NULL | 840.75 | NULL | NULL |
| Asparagus | NULL | NULL | 907.50 | NULL |

Result 7 ×                                                        ❶ Read Only

Output

Action Output

| # | Time | Action | Message |
|---|---|---|---|
| ✔ | 220 | 23:06:19 | WITH sold AS (select p.PRODUCT_NAME,t.Week_no ,sum(s.QUANTITY),ROW_... | 265 row(s) returned |
| ✔ | 221 | 23:11:04 | select p.PRODUCT_NAME, CASE t.QUARTER WHEN 1 THEN sum(s.TOTAL_SA... | 396 row(s) returned |

**Q5 Extract total sales of each product for the first and second half of year 2016 along with its total yearly sales.**

```sql
1 ● select p.PRODUCT_NAME,t.Mid,sum(s.TOTAL_SALE)
2    from data_warehouse_metro.sales s join data_warehouse_metro.time t
3    on(t.TIME_ID=s.TIME_ID)
4    join data_warehouse_metro.product p
5    on(p.PRODUCT_ID=s.PRODUCT_ID)
6    GROUP BY p.PRODUCT_NAME,t.Mid with rollup
7    order by p.PRODUCT_NAME,t.Mid;
8
```

**Result Grid** | Filter Rows: | Export: | Wrap Cell Content: 

| PRODUCT_NAME | Mid | sum(s.TOTAL_SALE) |
|---|---|---|
| NULL | NULL | 717674.19 |
| Apples | NULL | 4548.48 |
| Apples | 1 | 2274.24 |
| Apples | 2 | 2274.24 |
| Applesauce | NULL | 9863.80 |
| Applesauce | 1 | 4584.35 |
| Applesauce | 2 | 5279.45 |
| Asparagus | NULL | 3320.25 |

Result 2 ×                                                                    ⓘ Read Only

Output

Action Output ▾

| # | Time | Action | Message |
|---|---|---|---|
| ✓ | 223 | 23:17:47 | select p.PRODUCT_NAME,t.Mid,sum(s.TOTAL_SALE) from data_warehouse_metr... | 298 row(s) returned |
| ✓ | 224 | 23:17:51 | select p.PRODUCT_NAME,t.Mid,sum(s.TOTAL_SALE) from data_warehouse_metr... | 298 row(s) returned |

**Q6 Find an anomaly in the data warehouse dataset. write a query to show the anomaly and explain the anomaly in your project report.**

**Q7 Create a materialised view with the name "STOREANALYSIS_MV" that presents the product- wise sales analysis for each store.**

Limit to 50000 rows ▾

```
29 ●   CREATE VIEW STOREANALYSIS AS
30     SELECT s.store_id STORE_ID, p.product_id PROD_ID, SUM(l.total_sale) STORE_TOTAL
31     from data_warehouse_metro.store s, data_warehouse_metro.sales l,data_warehouse_metro.product p
32     where s.store_id = l.store_id and l.product_id = p.product_id
33      GROUP by s.store_name, p.product_name
34      ORDER by store_name, product_name;
35
36 ●   select * from data_warehouse_metro.STOREANALYSIS ;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: 𝘐A

| STORE_ID | PROD_ID | STORE_TOTAL |
|----------|---------|-------------|
| S-2 | P-1015 | 581.44 |
| S-2 | P-1080 | 1688.10 |
| S-2 | P-1000 | 527.25 |
| S-2 | P-1016 | 717.66 |
| S-2 | P-1030 | 402.93 |
| S-2 | P-1081 | 444.54 |
| S-2 | P-1017 | 681.05 |
| S-2 | P-1090 | 377.10 |

STOREANALYSIS 5 ✕     ❶ Read

Output

Action Output ▾

| # | Time | Action | Message |
|---|------|--------|---------|
| ✅ | 237 | 23:30:24 | CREATE VIEW STOREANALYSIS AS SELECT s.store_id STORE_ID, p.product_i... | 0 row(s) affected |
| ✅ | 238 | 23:30:24 | select * from data_warehouse_metro.STOREANALYSIS LIMIT 0, 50000 | 987 row(s) returned |