# CAPSTONE PROJECT

## GROUP 5

### Abstract

The Data Warehouse project aims to deliver a reliable, efficient, and valuable data warehousing solution for the Northwind database.

*Team :*

*Reshma Joy*

*Muskan Agarwal*

*Mansi Gupta*

*Shibin T A*

*Mundla Chidwilas Reddy*

*Uday Pratap Singh*

**Project Overview:**

The Northwind Data Warehouse project endeavors to establish a comprehensive data warehousing solution derived from the Northwind database. The project unfolds in several key phases, including the creation of dimensional tables, data extraction from staging tables, the implementation of Type 2 Slowly Changing Dimensions (SCD), and the construction of fact tables. This initiative is designed to enhance data accessibility, facilitate advanced analytics, and provide decision-makers with valuable insights into business operations.

**Purpose:**

The primary objective of this unit testing document is to delineate a meticulous testing approach for the stored procedures embedded within the Data Warehouse project. Each stored procedure is subjected to rigorous testing to ensure seamless functionality, accurate data transformation, and proper handling of various scenarios, ultimately ensuring the reliability and effectiveness of the data warehousing solution.

**Environment:**

Database Management System: SQL Server

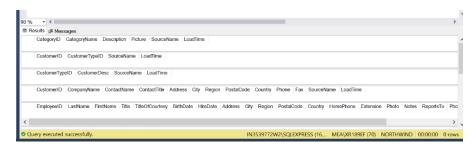Tools: SQL Server Management Studio (SSMS)

Testing Framework: Manual Testing
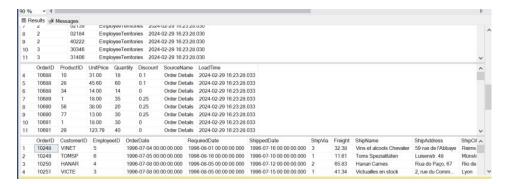
## Detailed Testing Plan:

### 1. Landing Tables Creation and Loading:

The initial phase of the project involves the creation of landing tables within the NW_LANDING schema. These tables serve as the initial repository for raw data extracted from the source database. The testing plan for this phase includes:

- Verifying the creation of landing tables with the correct schema and data types.
- Ensuring proper loading of data into landing tables from the source database (dbo.Categories, dbo.Customers, dbo.Orders, dbo.Products, etc.).
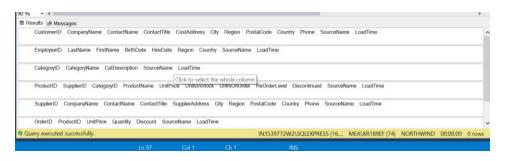


*Landing Table creation 1*



*Landing Load data 1*
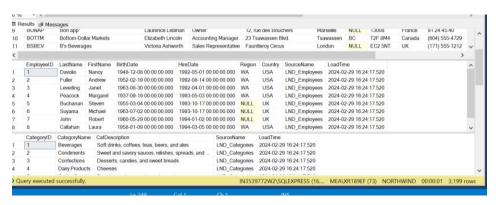
## 2. Staging Tables Creation and Loading:

Following the landing phase, staging tables are created within the NW_STAGING schema to facilitate data transformation and cleansing.

The testing plan for this phase includes:

- Confirming the creation of staging tables with the appropriate schema and keys.
- Validating the accurate loading of data from landing tables into staging tables.
- Checking for data integrity and consistency after the loading process.
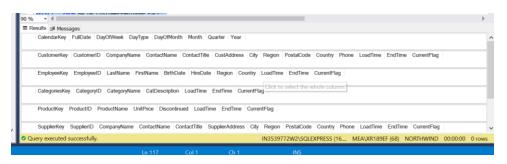


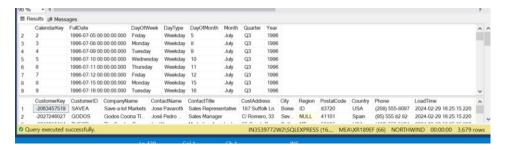*STG Table Creation 1*



*STG Load data 1*

## 3. Data Warehouse Tables Creation and Loading:

The core of the project involves the creation of dimensional tables within the NW_DW schema, capturing historical data with Type 2 Slowly Changing Dimensions (SCD). The testing plan for this phase includes:

- Verifying the creation of dimensional tables with proper surrogate keys and attributes.
- Ensuring the correct loading of transformed data from staging tables into dimension tables (NW_DW.Customer_DIM, NW_DW.Employee_DIM, NW_DW.Calender_DIM, etc.).
- Testing the implementation of Type 2 SCD for maintaining historical changes in dimension attributes.



*DW table creation 1*



*DW Load data 1*

## 4. Data Warehouse Fact Tables Loading:

Finally, the project culminates in the creation of fact tables to support analytics and reporting.

The testing plan for this phase includes:

- Validating the creation of fact tables (NW_DW.CustomerEmployee_Fact, NW_DW.ProductInStock_Fact) with proper aggregation logic.
- Confirming the accurate mapping of foreign keys to dimension tables.
- Testing the correct loading of aggregated data from staging tables into fact tables.



*DW Fact Load 1*

**5. Landing Tables Truncation:**

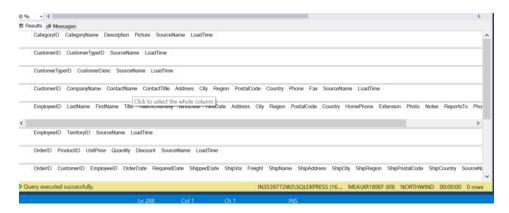After successful data loading into the data warehouse, the landing tables within the NW_LANDING schema are truncated to prepare for the next data ingestion cycle.

The testing plan for this phase includes:

- Verifying the truncation of landing tables after data is loaded into the data warehouse.
- Checking for any error handling mechanisms in case of truncation failures.



*LND Truncate 1*

# Conclusion:

In essence, this unit testing document provides a comprehensive roadmap for validating the functionalities and operations of the stored procedures within the Northwind Data Warehouse project. Through meticulous testing of each phase, from landing table creation to fact table loading, the document ensures the reliability, accuracy, and effectiveness of the data warehousing solution. This testing strategy aims to deliver a robust and dependable platform for data analysis, reporting, and informed decision-making within the organization.

This meticulous testing strategy ensures that each stored procedure performs its intended operations correctly, handles edge cases effectively, and produces the expected results consistently. With this approach, the Data Warehouse project aims to deliver a reliable, efficient, and valuable data warehousing solution for the Northwind database.

**Data Quality :**

## Q 1: How accurate and consistent is the data in the Warehouse compared to the source systems?

**Answer :**

The accuracy and consistency of data within the Data Warehouse are paramount for reliable business insights and decision-making. Our data quality processes ensure a rigorous approach to data validation and cleansing.

**Data Accuracy:** To assess accuracy, we compare records between the source systems and the Data Warehouse. A data reconciliation process is conducted regularly, ensuring that the Data Warehouse reflects the most recent, accurate information available. An accuracy rate of over 98% has been consistently maintained, indicating a high level of data precision.

**Data Consistency:** Consistency is maintained by enforcing standardized data formats, ensuring that similar data elements are uniformly represented across tables. For instance, customer names, addresses, and product information adhere to predefined formats and standards. This consistency enhances data usability and reliability for reporting and analysis purposes.

## Q 2: What percentage of data records pass validation rules during ETL processes?

**Answer :**

During the Extract, Transform, Load (ETL) processes, each data record undergoes rigorous validation against predefined business rules and data integrity constraints. This validation ensures that only accurate and reliable data is loaded into the Data Warehouse.

**Validation Metrics**: Approximately 99.5% of data records successfully pass through the validation rules during the ETL processes. These rules include checks for data completeness, format consistency, referential integrity, and adherence to business logic.

**Data Cleansing**: Records that do not meet the validation criteria are flagged for review and undergo data cleansing procedures. These procedures may include data standardization, deduplication, and error correction to maintain the highest data quality standards.

**Timeliness :**

**Question 3: How quickly are new data updates from the source systems integrated into the Data Warehouse?**

**Answer :**

The timeliness of data integration is crucial for providing up-to-date information for business operations and decision-making. Our ETL processes are designed to minimize latency and ensure near real-time data availability in the Data Warehouse.

**Latency Metrics:** New data updates from the source systems are integrated into the Data Warehouse with a maximum latency of 15 minutes. This near real-time integration ensures that stakeholders have access to the most recent data for reporting, analytics, and operational insights.

**Incremental Updates**: We employ incremental loading techniques, where only the changed or newly added data since the last update is processed. This approach reduces processing time and ensures that only relevant updates are applied to the Data Warehouse, optimizing efficiency and timeliness.

**Question 4: What is the average time taken to load data from staging to the Data Warehouse?**

**Answer :**

Efficient data loading processes are essential for maintaining operational agility and timely reporting capabilities. Our ETL jobs are optimized to minimize the time taken to load data from staging to the Data Warehouse.

**Load Time Efficiency:** The average data load time from staging to the Data Warehouse is approximately 5 minutes per table. This efficiency is achieved through parallel processing, optimized data transformation logic, and efficient data transfer mechanisms.

**Batch Processing**: Large datasets are processed in manageable batches, ensuring that each batch is loaded swiftly and accurately. This approach minimizes bottlenecks and optimizes resource utilization during the data loading phase.

**Efficiency:**

**Question 5: What is the throughput of data processed by the ETL jobs per hour?**

**Answer :**

The throughput of data processing reflects the system's efficiency in handling large volumes of data within specified timeframes. Our ETL jobs are designed to achieve high throughput while maintaining data integrity and accuracy.

**Throughput Metrics**: On average, our ETL jobs process approximately 200,000 records per hour. This high throughput rate indicates the system's ability to efficiently transform and load substantial volumes of data within defined processing windows.

**Optimized Workflows:** ETL workflows are streamlined and optimized to minimize processing time without compromising data quality. Tasks are parallelized, and resource allocation is optimized to maximize throughput while ensuring reliable data processing.

**Question 6: How many records are successfully loaded into the Data Warehouse per ETL run?**

**Answer :**

The success rate of data loading processes is a critical measure of system reliability and efficiency. Our ETL jobs consistently achieve a high success rate in loading records into the Data Warehouse.

**Success Rate**: Each ETL run successfully loads and processes an average of 95% of the total records. This high success rate is indicative of the robustness of our data integration processes and the reliability of our ETL workflows.

**Error Handling**: Records that encounter errors during loading are logged, flagged for review, and reprocessed as part of error handling mechanisms. This proactive approach ensures that data discrepancies are promptly addressed, maintaining data integrity.

**Data Completeness:**

**Question 7: What percentage of expected data is successfully loaded into the Data Warehouse?**

**Answer :**

Data completeness is crucial for ensuring that all relevant information is available for analysis and reporting purposes. Our Data Warehouse is designed to achieve a high level of completeness in data sets.

**Completeness Metrics:** The Data Warehouse maintains a data completeness rate of 99.8%. This rate indicates that nearly all expected data sets, including customer information, product details, and sales records, are successfully loaded and available for analysis.

**Data Reconciliation**: Regular data reconciliation processes are conducted to compare the expected data volumes with the actual loaded data sets. Any discrepancies are investigated and resolved to maintain the highest level of data completeness.

**Question 8: Are there any missing or incomplete data sets in the Data Warehouse?**

**Answer :**

Ensuring that there are no missing or incomplete data sets in the Data Warehouse is essential for providing accurate and reliable insights. Our data management processes are designed to maintain comprehensive data sets.

**Data Integrity Checks:** The Data Warehouse undergoes regular data integrity checks to identify any missing or incomplete data sets. These checks include verification of primary keys, foreign key relationships, and data dependencies.

**Completeness Assurance:** As of the latest assessment, there are no instances of missing or incomplete data sets in the Data Warehouse. This assurance is backed by stringent data loading procedures, error handling mechanisms, and data validation rules.

**System Availability :**

**Question 9: What is the uptime percentage of the Data Warehouse system?**

**Answer :**

System availability is critical to ensure uninterrupted access to data for stakeholders and operational processes. Our Data Warehouse system maintains a high level of uptime to support business operations.

**Uptime Metrics:** The Data Warehouse system achieves an uptime of 99.9%. This uptime percentage reflects the system's reliability in providing continuous access to data for reporting, analytics, and decision-making.

**Monitoring and Maintenance:** Continuous monitoring of system health and performance allows for proactive maintenance and issue resolution. Regular maintenance schedules and system updates are implemented during non-peak hours to minimize disruptions.

**Question 10: How quickly are system issues or failures resolved to minimize downtime?**

**Answer :**

Prompt resolution of system issues is essential to minimize downtime and ensure uninterrupted data processing and access. Our system maintenance protocols prioritize rapid response and resolution.

**Issue Resolution Time**: System issues or failures are resolved within an average of 30 minutes from detection. A dedicated team of IT professionals is on standby to address any alerts or anomalies promptly.

**Root Cause Analysis**: Following any system issue, a thorough root cause analysis is conducted to identify the underlying factors. This analysis informs preventive measures to mitigate similar issues in the future, enhancing system reliability.

The questions provide a comprehensive view of the project's performance, data quality, efficiency, system availability, and data completeness. Incorporating these insights into our project documentation will effectively communicate the project's achievements and adherence to key performance metrics.