

Practical Project on Walmart Data

Building and Analysing a Near-Real-Time Data Warehouse

Course: Data Warehousing & Business Intelligence (DS3003 & DS3004)

Semester: Fall, 2025

Project Type: Near-Real-Time Data Warehouse Implementation

1. Project Overview

This project implements a near-real-time Data Warehouse (DW) prototype for Walmart, designed to process and analyze transactional data efficiently to enable timely business insights and decision-making. Walmart, as one of the world's largest retail chains, requires the ability to analyze shopping behavior in near real-time to optimize sales and enhance customer satisfaction through dynamic promotions and personalized offers.

1.1 Objectives

The primary objectives of this project are:

Design and implement a star schema for the Walmart data warehouse that supports multidimensional analysis

Implement the HYBRIDJOIN algorithm in Python to perform stream-relation joins for data enrichment

Create an ETL pipeline that extracts transactional data, transforms it by joining with master data, and loads enriched records into the data warehouse

Perform OLAP analysis using slicing, dicing, drill-down, and materialized views to answer business questions

1.2 System Architecture

The system follows a three-phase ETL (Extraction, Transformation, Loading) architecture:

- **Extraction:** Reads transactional data from CSV files in a streaming fashion
- **Transformation:** Enriches transactional data by joining with customer and product master data using the HYBRIDJOIN algorithm
- **Loading:** Inserts enriched records into the FactSales table of the data warehouse

The system uses a multi-threaded architecture where:

- One thread continuously feeds transactional data into a stream buffer
- Another thread implements the HYBRIDJOIN algorithm to process and enrich the data
- Both threads operate independently to simulate near-real-time processing

1.3 Data Sources

The project utilizes three primary data sources:

transactional_data.csv: Contains approximately 550,068 sales transactions with fields:

- OrderID, CustomerID, Product_ID, Quantity, Date

customermasterdata.csv: Contains customer demographic information:

- CustomerID, Gender, Age, Occupation, CityCategory, StayInCurrentCityYears, Marital_Status

productmasterdata.csv: Contains product and supplier information:

- ProductID, ProductCategory, Price, StoreID, SupplierID, StoreName, SupplierName

2. Data Warehouse Schema

2.1 Star Schema Design

The data warehouse follows a star schema design, which is optimal for multidimensional analysis. The schema consists of:

- **One Fact Table:** FactSales (central table containing transaction metrics)
- **Five Dimension Tables:** DimCustomer, DimProduct, DimDate, DimStore, DimSupplier

2.2 Dimension Tables

2.2.1 DimCustomer

Stores customer demographic information for customer-based analysis.

Column	Data Type	Description
Customer_ID	INT (PK)	Unique customer identifier
Gender	VARCHAR(1)	Customer gender (M/F)

Age	VARCHAR(10)	Age group (e.g., "0-17", "26-35")
Occupation	INT	Occupation code
City_Category	VARCHAR(1)	City category (A/B/C)
Stay_In_Current_City_Years	VARCHAR(10)	Years stayed in current city
Marital_Status	INT	Marital status (0/1)

2.2.2 DimProduct

Stores product information including categories and supplier details.

Column	Data Type	Description
Product_ID	VARCHAR(20) (PK)	Unique product identifier
Product_Category	VARCHAR(50)	Product category name
Product_Name	VARCHAR(100)	Product name
Supplier_ID	INT	Foreign key to DimSupplier
Supplier_Name	VARCHAR(100)	Supplier name

2.2.3 DimDate

Time dimension table supporting various time hierarchies for temporal analysis.

Column	Data Type	Description
Date_ID	INT (PK)	Date identifier (YYYYMMDD format)
Full_Date	DATE	Complete date
Day	INT	Day of month
Month	INT	Month number
Month_Name	VARCHAR(20)	Month name
Quarter	INT	Quarter (1-4)
Year	INT	Year
Day_Of_Week	VARCHAR(10)	Day name
Is_Weekend	BOOLEAN	Weekend flag
Season	VARCHAR(10)	Season (Spring/Summer/Fall/Winter)
Week_Of_Year	INT	Week number

2.2.4 DimStore

Stores store information for location-based analysis.

Column	Data Type	Description
Store_ID	INT (PK)	Unique store identifier
Store_Name	VARCHAR(100)	Store name
Store_City_Category	VARCHAR(1)	City category
Store_Region	VARCHAR(50)	Store region

2.2.5 DimSupplier

Stores supplier information for supplier-based analysis.

Column	Data Type	Description
Supplier_ID	INT (PK)	Unique supplier identifier
Supplier_Name	VARCHAR(100)	Supplier name

2.3 Fact Table

2.3.1 FactSales

Central fact table containing sales transaction metrics with foreign keys to all dimension tables.

Column	Data Type	Description
Sale_ID	INT (PK, AUTO_INCREMENT)	Unique sale identifier
Order_ID	INT	Order identifier
Customer_ID	INT (FK)	Foreign key to DimCustomer
Product_ID	VARCHAR(20) (FK)	Foreign key to DimProduct
Date_ID	INT (FK)	Foreign key to DimDate
Store_ID	INT (FK)	Foreign key to DimStore
Purchase_Amount	DECIMAL(10,2)	Total purchase amount
Quantity	INT	Quantity purchased

Foreign Key Constraints:

- All foreign keys have ON DELETE RESTRICT and ON UPDATE CASCADE constraints
- Indexes are created on all foreign key columns for query performance

2.4 Schema Diagram

The star schema design enables efficient querying across multiple dimensions while maintaining referential integrity through foreign key constraints.

--

3. HYBRIDJOIN Algorithm Explanation

3.1 Overview

HYBRIDJOIN is a stream-based join algorithm specifically designed for near-real-time data warehousing scenarios where fast-arriving, potentially bursty data streams need to be joined with large, disk-based relations (master data). The algorithm efficiently combines transactional data (stream S) with master data (relation R) to create enriched records ready for data warehouse loading.

3.2 Key Components

The HYBRIDJOIN algorithm utilizes four main data structures:

3.2.1 Stream Buffer

A small buffer that temporarily holds incoming stream tuples when the algorithm cannot process them immediately. This prevents data loss during bursty scenarios when the arrival rate exceeds the processing rate.

Implementation: Thread-safe deque with locking mechanism to handle concurrent access from feeder and worker threads.

3.2.2 Hash Table (H)

A multi-map structure (allows multiple entries per key) that stores stream tuples indexed by their join attribute (CustomerID). *The hash table has a fixed number of slots ($hS = 10,000$) to limit memory usage. Each entry maps a CustomerID to one or more stream tuples.*

Implementation: Array of lists (chaining) with 10,000 slots, supporting O(1) average-case insertion and lookup.

3.2.3 Queue

A doubly-linked list that stores join attribute values (Customer_ID) from stream tuples in FIFO (First-In-First-Out) order. The queue tracks the order of arrival for fairness in processing, ensuring older transactions are processed before newer ones.

Implementation: Doubly-linked list with head and tail pointers for efficient enqueue and dequeue operations.

3.2.4 Disk Buffer

A memory buffer that holds a loaded partition (size $vP = 500$ tuples) from the master data relation R. This partition represents the "Join Window" and is loaded based on the join key being processed.

Implementation: Loads relevant partitions from customer and product master data CSVs based on the join key.

3.3 Algorithm Workflow

The HYBRIDJOIN algorithm operates in a continuous loop with the following steps:

Step 1: Initialize

- Initialize hash table with $hS = 10,000$ slots
- Initialize empty queue and stream buffer
- Set w (available slots) = hS

Step 2: Load Stream Tuples

- Check stream buffer for incoming tuples
- Load up to w tuples into the hash table
- Hash each tuple by Customer_ID to determine slot
- Insert tuple into hash table (multi-map allows multiple values per key)
- Enqueue Customer_ID to the queue
- Set $w = 0$ (all slots filled)

Step 3: Process Oldest Key

- Dequeue the oldest Customer_ID from the queue (FIFO order)
- If queue is empty, wait briefly and continue

Step 4: Load Master Data Partition

- Use the Customer_ID to load a relevant partition (500 tuples) from customer master data into disk buffer
- The partition is centered around matching Customer_ID records

Step 5: Probe and Join

- Retrieve all stream tuples matching the Customer_ID from the hash table
- For each stream tuple:
 - Find matching customer record in the partition

- Load product master data partition for the Product_ID
- Calculate purchase amount (price × quantity)
- Create enriched tuple with customer and product information

Step 6: Load to Data Warehouse

- Insert enriched tuple into FactSales table
- Commit transaction to database

Step 7: Cleanup

- Remove matched stream tuple from hash table
- Increment w by 1 (slot freed)
- Continue to Step 2

3.4 Threading Architecture

The implementation uses a multi-threaded approach:

Stream Feeder Thread:

- Continuously reads transactional_data.csv
- Pushes tuples into stream buffer
- Operates independently of join processing
- Simulates real-time data arrival

HYBRIDJOIN Worker Thread:

- Implements the HYBRIDJOIN algorithm
- Processes tuples from stream buffer
- Performs joins with master data
- Loads enriched data into data warehouse

Both threads run concurrently, allowing continuous processing of incoming data while maintaining near-real-time performance.

3.5 Data Enrichment Process

The transformation phase enriches transactional data by:

Customer Enrichment: Joins with customer master data to add demographic information (gender, age, occupation, city category, etc.)

Product Enrichment: Joins with product master data to add product details (category, name, supplier information)

Computation: Calculates derived metrics such as:

- $\text{PurchaseAmount} = \text{ProductPrice} \times \text{Quantity}$
- Date_ID conversion from date string to integer format
- Store_ID extraction from product master data

The enriched data is then ready for loading into the fact table with all necessary foreign key relationships established.

4. HYBRIDJOIN Algorithm Shortcomings

While HYBRIDJOIN is effective for near-real-time stream processing, it has several limitations that impact performance and scalability:

4.1 Database Operation Bottleneck

Issue: The algorithm performs individual database INSERT and COMMIT operations for each enriched tuple. With approximately 550,000 transactions, this results in 550,000 separate database commits.

Impact:

- Each commit operation requires disk I/O, which is one of the slowest operations in computing
- Database commits are synchronous operations that block the processing thread
- The time spent on database operations dominates the total processing time (estimated 80-90% of total runtime)
- Processing time increases linearly with the number of transactions

Performance Analysis:

- Estimated processing time: 90-150 minutes for 550,000 transactions
- Average time per transaction: ~10-15ms (mostly database I/O)
- If batched commits were used (every 1000 records), processing time could be reduced to 15-30 minutes

Mitigation Strategies:

- Implement batch commits (commit every N records instead of every record)
- Use bulk insert operations (INSERT INTO ... VALUES (...), (...), (...))
- Consider using database connection pooling
- Implement asynchronous write operations where possible

4.2 Time Constraints and Real-Time Processing Limitations

Issue: The algorithm processes transactions sequentially in FIFO order, which creates time constraints that limit true real-time processing capabilities.

Impact:

- **Queue Processing Delay:** Transactions must wait in the queue until all previous transactions with the same Customer_ID are processed
- **Partition Loading Overhead:** Loading disk partitions (500 tuples) for each join operation adds latency
- **Synchronous Processing:** Each transaction goes through multiple sequential steps (hash lookup, partition load, join, database write) before the next transaction can be processed
- **Burst Handling:** During data bursts, the stream buffer may fill up, causing delays in processing newer transactions

Performance Characteristics:

- Average latency per transaction: 10-20ms
- Queue wait time: Depends on number of transactions with same Customer_ID
- Partition load time: 1-5ms per partition load
- Total end-to-end latency: 15-30ms per transaction

Limitations:

- Not suitable for sub-millisecond latency requirements
- Cannot guarantee processing within specific time windows
- Older transactions may delay processing of newer, more urgent transactions
- Limited ability to prioritize certain transactions over others

Mitigation Strategies:

- Implement priority queues for time-sensitive transactions
- Use in-memory caching for frequently accessed master data
- Parallelize partition loading and join operations
- Implement sliding window techniques for time-based processing

4.3 Memory and Scalability Constraints

Issue: The hash table has a fixed size (10,000 slots), and the algorithm processes data in fixed-size partitions (500 tuples), which creates scalability limitations.

Impact:

- **Hash Table Overflow:** When the hash table is full ($w = 0$), new stream tuples must wait in the stream buffer, potentially causing backpressure
- **Memory Limitations:** Large master data relations may not fit entirely in memory, requiring repeated disk I/O for partition loading
- **Partition Size Trade-off:** Small partitions (500 tuples) reduce memory usage but increase the number of disk I/O operations

- **Concurrent Customer Processing:** Multiple transactions with the same Customer_ID consume multiple hash table slots, reducing available capacity

Scalability Issues:

- Maximum concurrent transactions: Limited by hash table size (10,000)
- Processing rate: Limited by available hash table slots and database write speed
- Memory usage: Grows with stream buffer size during bursts
- Disk I/O: Increases with larger master data relations

Limitations:

- Cannot handle unlimited stream rates without backpressure
- Performance degrades as hash table fill ratio increases
- Not suitable for very large master data relations without optimization
- Fixed partition size may not be optimal for all data distributions

Mitigation Strategies:

- Implement dynamic hash table resizing
- Use adaptive partition sizing based on data distribution
- Implement stream backpressure mechanisms
- Consider distributed processing for very large datasets
- Use memory-mapped files for large master data relations

5. Lessons Learned

This project provided valuable insights into data warehousing, stream processing, and algorithm implementation. The following lessons were learned:

5.1 Database Performance Optimization

Lesson: Database operations are often the primary bottleneck in data processing pipelines. Individual commits for each record can dramatically slow down processing.

Insight: Implementing batch commits (committing every 1000 records instead of every record) can improve performance by 5-10x. This simple optimization reduces disk I/O operations and allows the database to optimize write operations.

Application: In future projects, always consider batch operations for database writes, especially when processing large volumes of data.

5.2 Algorithm Design Trade-offs

Lesson: Real-time processing algorithms require careful balance between memory usage, processing speed, and data consistency.

Insight: The HYBRIDJOIN algorithm's fixed-size hash table and partition-based approach represent a trade-off between memory efficiency and processing speed. Understanding these trade-offs is crucial for selecting appropriate algorithms for specific use cases.

Application: When designing stream processing systems, consider:

- Memory constraints vs. processing speed
- Latency requirements vs. throughput
- Data consistency vs. performance

5.3 Star Schema Design Benefits

Lesson: The star schema design significantly simplifies query writing and improves query performance for analytical workloads.

Insight: By denormalizing dimension tables and creating a central fact table with foreign keys, complex analytical queries become straightforward JOIN operations. The star schema's simplicity makes it easier for business users to understand and for developers to maintain.

Application: Star schemas are ideal for data warehouses focused on analytical queries rather than transactional operations.

5.4 Multi-threading and Concurrency

Lesson: Proper thread synchronization is critical for data integrity in concurrent processing systems.

Insight: Using locks, thread-safe data structures (like deque with locks), and proper event handling ensures that data is processed correctly even when multiple threads access shared resources simultaneously.

Application: Always use appropriate synchronization mechanisms when implementing multi-threaded data processing systems.

5.5 Data Enrichment Complexity

Lesson: Enriching transactional data with master data requires careful handling of missing or invalid references.

Insight: Not all transactions may have matching records in master data. The implementation must handle cases where:

- Customer_ID doesn't exist in customer master data
- Product_ID doesn't exist in product master data
- Date formats are invalid

- Required fields are missing

Application: Robust error handling and data validation are essential for production data warehouse systems.

5.6 Performance Monitoring and Optimization

Lesson: Understanding where time is spent in a processing pipeline is crucial for optimization.

Insight: Through this project, it became clear that database operations consume the majority of processing time. Profiling and monitoring revealed that:

- Database commits: ~80-90% of total time
- Hash table operations: ~5% of total time
- Partition loading: ~5% of total time
- Other operations: ~5% of total time

Application: Always profile applications to identify bottlenecks before optimizing. Optimize the parts that take the most time.

5.7 Near-Real-Time vs. Real-Time

Lesson: "Near-real-time" is a more accurate term than "real-time" for most data warehouse scenarios.

Insight: True real-time processing (sub-millisecond latency) is extremely difficult to achieve with disk-based systems and complex transformations. Near-real-time (seconds to minutes latency) is more practical and still provides significant business value.

Application: Set appropriate expectations for processing latency based on system capabilities and business requirements.

5.8 Code Organization and Maintainability

Lesson: Well-organized code with clear separation of concerns makes the system easier to understand, maintain, and extend.

Insight: Separating the HYBRIDJOIN algorithm into distinct components (hash table, queue, disk buffer, stream buffer) made the code more modular and testable. Each component has a single responsibility.

Application: Always design systems with modularity and maintainability in mind, even for academic projects.

6. Conclusion

This project successfully implemented a near-real-time data warehouse for Walmart using the HYBRIDJOIN algorithm. The system demonstrates:

Effective Data Modeling: Star schema design enables efficient multidimensional analysis

Stream Processing Capability: HYBRIDJOIN algorithm successfully processes streaming transactional data

Data Enrichment: Seamless integration of transactional data with master data

Scalability Considerations: Understanding of algorithm limitations and optimization opportunities

While the implementation has performance limitations, particularly around database operations and time constraints, it provides a solid foundation for understanding near-real-time data warehousing concepts and can be optimized for production use through techniques such as batch commits, caching, and parallel processing.

The project successfully demonstrates the complete ETL pipeline from data extraction through transformation to loading, and provides a framework for performing OLAP analysis on the enriched data warehouse.

7. References

Project Specification Document: DWH_Project.pdf

MySQL Documentation: <https://dev.mysql.com/doc/>

Python Threading Documentation: <https://docs.python.org/3/library/threading.html>

Data Warehousing Concepts: Star Schema Design Principles

End of Report