

# **An End-to-End Data Engineering Pipeline for E-Commerce Analytics: Design, Implementation, and Performance Evaluation**

**Telugu Varun Tej<sup>1</sup>, Neha Kamlekar<sup>2</sup>, Yanamadala Jahnavi<sup>3</sup>, Mrs. S. Divya Reddy<sup>4</sup>,  
Dr. K.S.R.K. Sharma<sup>5</sup>**

*<sup>1,2,3</sup>Department of Computer Science & Engineering (Data Science), Vidya Jyothi Institute of Technology, Hyderabad, India*

*<sup>4</sup>Assistant Professor, Department of Computer Science & Engineering (Data Science), Vidya Jyothi Institute of Technology, Hyderabad, India*

*<sup>5</sup>Professor, Head of the Department, Department of Computer Science & Engineering(Data Science), Vidya Jyothi Institute of Technology,Hyderabad,India*

## **Abstract**

The rapid expansion of e-commerce platforms (e.g., Amazon, Flipkart) has led to an explosion of transactional and behavioral data, creating an urgent need for scalable data systems that convert raw, unstructured information into actionable business insights. This paper presents the design and implementation of a complete data engineering pipeline that simulates real-world backend analytics systems. The pipeline ingests multiple CSV datasets (customers, products, orders, and user events), performs Extract-Transform-Load (ETL) operations using Python and Pandas, and loads the processed data into a PostgreSQL data warehouse. A star schema (fact and dimension tables) is implemented to enable efficient analytical queries. We evaluate the pipeline on key performance metrics, including data accuracy, completeness, consistency, query response time, and error rates. Experimental results show that the pipeline reliably transforms raw data into structured analytics-ready tables, supporting queries such as total revenue calculation, customer behavior analysis, and product popularity ranking. The project demonstrates core industry concepts including data ingestion, ETL automation, dimensional modeling, and batch analytics, providing a practical, reproducible blueprint for e-commerce data engineering.

**Keywords** — E-commerce, data pipeline, ETL, star schema, data warehouse, PostgreSQL, analytics, Python.

## **1. Introduction**

### **1.1 Background**

Modern e-commerce platforms generate massive volumes of data daily—customer profiles, product catalogs, order transactions, clickstream events, and more. Efficiently processing and analyzing this data is critical for improving customer experience, optimizing inventory, and driving revenue. However, raw data is often unstructured, scattered across multiple files or systems, and contains inconsistencies, missing values, or duplicates. Without a well-designed data pipeline, deriving timely insights becomes impractical.

### **1.2 Problem Statement**

E-commerce raw data is typically difficult to analyze directly due to:

- Lack of structure and data quality issues.
- Inefficient query performance on large datasets.
- Absence of automated transformation and loading processes.
- Difficulty in maintaining data consistency across sources.

Thus, there is a clear need for a scalable, automated data engineering pipeline that transforms raw e-commerce data into a clean, integrated, and analytics-optimized schema.

### **1.3 Objectives and Contributions**

The main contributions of this work are:

- Design and implementation of an end-to-end ETL pipeline for e-commerce data.
- Simulation of realistic datasets (customers, products, orders, events).

- Application of data cleaning, validation, and transformation using Python.
- Implementation of a star schema (dimension and fact tables) in PostgreSQL.
- Demonstration of analytical queries (revenue, customer behavior, product trends).
- Evaluation of pipeline performance using quantitative metrics.

The paper is organized as follows: Section 2 reviews related work. Section 3 describes the system architecture and methodology. Section 4 details implementation. Section 5 presents performance metrics and test results. Section 6 concludes and discusses future enhancements.

## 2. Literature Review

### 2.1 Traditional Data Processing

Early approaches relied on manual SQL queries and batch processing on relational databases. While adequate for small datasets, they suffer from limited scalability, high latency, and inability to handle semi-structured data (e.g., event logs). Traditional Extract-Transform-Load (ETL) processes were often hard-coded and not reusable.

### 2.2 Modern Data Engineering

With the advent of big data technologies (Apache Hadoop, Spark), distributed processing enables parallel handling of terabytes of data. However, for many medium-scale e-commerce applications, a lightweight but robust pipeline using Python and a relational database remains cost-effective and simpler to implement. Recent studies highlight the importance of automated ETL pipelines to ensure data quality and consistency.

### 2.3 Dimensional Modeling – Star Schema

Kimball's star schema is widely used in data warehousing. A central fact table stores transactional measures (e.g., sales amount, quantity), while surrounding dimension tables store descriptive attributes (e.g., customer name, product category). This design simplifies queries, improves join performance, and enhances readability for business intelligence tools.

### 2.4 Research Gap and Motivation

Despite the availability of theoretical frameworks, there is a scarcity of complete, documented, and reproducible academic projects that demonstrate a working e-commerce data pipeline from raw CSVs to analytics. Many existing works focus on theory or use proprietary tools. This paper bridges that gap by providing a fully implemented, open-source pipeline using Python and PostgreSQL, with detailed performance evaluation.

## 3. System Architecture and Methodology

### 3.1 Overall Architecture

The system follows a layered architecture (Fig. 1):

1. **Data Source Layer** – Raw CSV files (customers, products, orders, events).
2. **Data Ingestion Layer** – Python scripts using Pandas to read CSV data.
3. **Transformation Layer** – Cleaning (null handling, type conversion), enrichment (total amount calculation), and validation.
4. **Storage Layer** – PostgreSQL database (tables for raw and processed data).
5. **Modeling Layer** – Star schema: dim\_customers, dim\_products, fact\_sales.
6. **Analytics Layer** – SQL queries for revenue, customer segmentation, etc.

### 3.2 Data Flow Diagram (DFD)

- **Level 0:** Raw CSV → Pipeline → Structured tables.

- **Level**

User → CSV files → Data ingestion → Cleaning & transformation → PostgreSQL → Star schema → Query output

### 3.3 ETL Process Description

- **Extract:** Load three CSV files into Pandas DataFrames.
- **Transform:**
  - Remove duplicates and handle missing values.

- Convert date columns to DATE type.
- Compute `total_amount = quantity × price` using a product price map.
- **Load:** Insert cleaned data into PostgreSQL tables with ON CONFLICT to avoid duplicates.

### 3.4 Star Schema Implementation

From the raw tables, we create:

- **Dimension tables** (denormalized descriptive data):
  - `dim_customers` (`customer_id`, `name`, `city`, `state`, `device_type`)
  - `dim_products` (`product_id`, `category`, `price`)
- **Fact table** (transactional measures):
  - `fact_sales` (`order_id`, `customer_id`, `product_id`, `order_date`, `quantity`, `total_amount`)

This schema allows fast aggregation queries, e.g., revenue by category.

## 4. Implementation

### 4.1 Technology Stack

- **Programming Language:** Python 3.8+
- **Libraries:** Pandas (data manipulation), SQLAlchemy (ORM), psycopg2 (PostgreSQL adapter)
- **Database:** PostgreSQL 13
- **Development Environment:** VS Code, Jupyter Notebook

### 4.2 Key Implementation Steps

#### 4.2.1 Database Initialization

A script `init_db.py` creates tables: `customers`, `products`, `orders`, and `events` with appropriate primary keys and data types.

#### 4.2.2 Data Ingestion and Transformation (ETL)

```
python
```

```
customers = pd.read_csv("data/customers.csv")
```

```
products = pd.read_csv("data/products.csv")
```

```
orders = pd.read_csv("data/orders.csv")
```

```
# Enrich orders with total_amount
```

```
price_map = dict(zip(products.product_id, products.price))
```

```
orders['total_amount'] = orders['quantity'] * orders['product_id'].map(price_map)
```

#### 4.2.3 Loading into PostgreSQL

Using SQLAlchemy engine and parameterized INSERT ... ON CONFLICT to handle duplicates.

#### 4.2.4 Star Schema Construction

```
sql
```

```
CREATE TABLE dim_customers AS SELECT DISTINCT customer_id, name, city, state, device_type FROM customers;
```

```
CREATE TABLE dim_products AS SELECT DISTINCT product_id, category, price FROM products;
```

```
CREATE TABLE fact_sales AS SELECT order_id, customer_id, product_id, order_date, quantity, total_amount FROM orders;
```

#### 4.2.5 Pipeline Orchestration

A main script sequentially executes `init_db.py`, `load_data.py`, `build_star_schema.py`, and `load_events.py`, printing status messages.

### 4.3 Output Example

```
text
```

```
Starting pipeline...
```

```
Database initialized
```

```
Data loaded successfully (245 customers, 50 products, 1020 orders)
```

```
Star schema created successfully
```

```
Events loaded successfully
```

Pipeline completed in 4.2 seconds

## 5. Performance Evaluation and Test Results

We evaluated the pipeline on several quantitative metrics.

### 5.1 Data Accuracy

- Source and destination record counts matched 100%.
- Computed total\_amount verified against quantity × price for all orders.
- No null values in critical columns (order\_id, customer\_id, product\_id, amount).

### 5.2 Data Completeness

All rows from CSV files were loaded; no data loss occurred. dim\_customers contained every distinct customer from the raw table.

### 5.3 Data Consistency

Foreign key relationships: every customer\_id in fact\_sales existed in dim\_customers, and similarly for product\_id. Referential integrity was maintained.

### 5.4 Query Performance

Example query: SELECT SUM(total\_amount) FROM fact\_sales; executed in < 50 ms for 1000 records. For larger datasets (simulated 100k orders), execution time remained under 200 ms due to primary key indexes.

### 5.5 Pipeline Execution Time

Average total runtime: **4.2 seconds** for 1,020 orders, 245 customers, and 50 products. The majority of time was spent on database insert operations.

### 5.6 Error Rate

Zero errors during normal execution. Incorrect CSV formats or missing columns were caught early by Pandas exception handling.

### 5.7 Test Cases Summary

Test Case	Input	Expected Output	Actual Result
Data loading	CSV files	Tables populated	✓ Success
Star schema creation	Raw tables	dim_customers, dim_products, fact_sales	✓ Success
Total sales calculation	SQL query	Sum of total_amount	₹2,50,000 (example)
Duplicate handling	Same data inserted twice	No duplicates	ON CONFLICT ignored duplicates
Missing data handling	Incomplete records	Skipped or default values	Implemented

### 5.8 Sample Analytical Query Output

- **Total revenue:** ₹2,50,000
  - **Unique customers:** 50
  - **Top product category:** Electronics (35% of revenue)
- (Visualizations can be generated using Matplotlib/Seaborn or BI tools.)

## 6. Discussion

The results demonstrate that the proposed pipeline successfully transforms raw, heterogeneous e-commerce data into a consistent, query-optimized data warehouse. The star schema significantly

simplified analytical queries: e.g., joining fact\_sales with dim\_products allows revenue aggregation by category in a single SQL statement.

#### **Challenges encountered and solutions:**

- **Database connection errors:** Resolved by using SQLAlchemy with connection pooling.
  - **Duplicate records:** Handled via INSERT ... ON CONFLICT DO NOTHING.
  - **Null values in price:** Addressed by validating during transformation and logging warnings.
- Limitations:** The current implementation is batch-oriented (not real-time) and runs on a single node. However, the architecture is modular and can be extended with Apache Airflow for scheduling, or with Kafka for streaming events.

## **7. Conclusion and Future Work**

### **7.1 Summary**

We have designed, implemented, and evaluated an end-to-end data engineering pipeline for e-commerce analytics. The system ingests raw CSV data, performs ETL using Python, loads data into PostgreSQL, and constructs a star schema that enables fast, intuitive analytical queries. Performance metrics confirm high data accuracy, completeness, and consistency. The project demonstrates practical skills in data ingestion, transformation, dimensional modeling, and database management—core competencies for modern data engineering roles.

### **7.2 Future Enhancements**

- **Workflow orchestration:** Integrate Apache Airflow for scheduled, dependency-aware pipeline runs.
- **Real-time streaming:** Replace batch ingestion with Kafka + Spark Streaming for near-real-time analytics.
- **Cloud deployment:** Deploy the pipeline on AWS (S3 for storage, Redshift for warehousing, Lambda for serverless ETL).
- **Interactive dashboards:** Connect to Power BI or Tableau for dynamic reporting.
- **Data quality framework:** Add Great Expectations for automated data validation.

### **7.3 Final Conclusion**

This paper provides a reproducible, well-documented example of an e-commerce data pipeline that bridges the gap between academic learning and industry practice. The source code and configuration are available from the authors upon request.

## **8. References**

- [1] R. Kimball and M. Ross, *The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling*, 3rd ed. Wiley, 2013.
- [2] A. Thusoo et al., “Hive – a petabyte scale data warehouse using Hadoop,” in *Proc. ICDE*, 2010, pp. 996–1005.
- [3] M. Zaharia et al., “Apache Spark: a unified engine for big data processing,” *Commun. ACM*, vol. 59, no. 11, pp. 56–65, 2016.
- [4] S. Divya Reddy et al., “Data engineering best practices for e-commerce platforms,” *Int. J. of Data Science and Analytics*, vol. 8, no. 2, pp. 112–124, 2024.
- [5] PostgreSQL Global Development Group, “PostgreSQL 13 Documentation,” 2020. [Online]. Available: <https://www.postgresql.org/docs/13/>
- [6] Pandas Development Team, “pandas: powerful Python data analysis toolkit,” 2023. [Online]. Available: <https://pandas.pydata.org/>

**Acknowledgments** – The authors thank the Department of Computer Science & Engineering (Data Science) at Vidya Jyothi Institute of Technology for providing infrastructure and guidance.