

# LEGO Database ETL Pipeline Report

## Executive Summary

This report documents the implementation of a comprehensive Extract, Transform, Load (ETL) pipeline and data mart for the LEGO Rebrickable dataset. The solution encompasses a relational database architecture, data transformation processes, dimensional modeling techniques, and automated reporting procedures. Through careful consideration of database design principles and ETL methodologies, the implementation provides a great foundation for analyzing LEGO product data, supporting efficient querying capabilities and business intelligence functions.

## 1. Database Architecture and Implementation Strategy

The database architecture was implemented using two distinct methodologies, each with specific advantages and considerations that influenced the final design decisions.

### 1.1 Implementation Methodologies

#### Method 1: Dynamic Schema Generation with Post-Implementation Refinement

The first implementation approach allowed for dynamic table creation during the data loading process. This method exemplifies an agile database development paradigm where schema structures evolve based on actual data characteristics rather than predetermined specifications. Tables were automatically generated with default data types assigned according to the content of the imported CSV files.

This approach offers distinct advantages in exploratory data projects where schema requirements may not be fully known at the outset. It permits rapid initial deployment and provides flexibility to adapt to unexpected data characteristics. Following data insertion, a systematic schema refinement process was conducted to implement appropriate constraints, normalizing the database structure while preserving data integrity.

Post-loading schema modifications included:

- Altering null-allowing columns to not null where appropriate
- Adding primary key constraints to establish entity identification
- Implementing foreign key constraints to maintain referential integrity
- Creating composite keys for junction tables representing many-to-many relationships

This methodology demonstrates effective adaptation of database structures to accommodate data realities while maintaining relational database principles.

#### Method 2: Pre-defined Schema with Controlled Data Loading

The second implementation employed a more traditional database design approach with pre-defined schema structures. Tables were created with explicitly defined data types, constraints, and relationships prior to data loading. This methodology reflects a data engineering best practice where database architecture is fully specified before integration with source data.

The pre-defined schema approach offers several advantages:

- Ensures data quality through enforced constraints during the loading process
- Eliminates the need for post-loading structural modifications
- Provides clear documentation of the intended data model
- Reduces the risk of data type incompatibilities and truncation issues

This method required additional preparation but resulted in a more controlled data environment with guaranteed structural integrity from inception.

## **2. Data Transformation and Quality Management**

### **2.1 Data Preprocessing and Cleansing**

Data quality assessment formed a critical preliminary stage in the ETL pipeline. Initial analysis revealed several quality issues within the source CSV files that required resolution before database integration:

The data cleansing process employed statistical methods for handling missing values:

- Numeric field nulls were replaced with mean values to maintain statistical consistency
- Categorical field nulls were populated with mode values to reflect the most common occurrences
- Foreign key null references were systematically addressed to maintain referential integrity

This methodical approach to null value handling ensured complete datasets while minimizing statistical distortion, enabling

successful database loading without compromise to analytical capabilities.

### **2.2 ETL Process Implementation**

The ETL pipeline was implemented using Python, leveraging libraries such as Pandas for data manipulation and SQLAlchemy for database interaction. The process incorporated these key components:

Data extraction involved reading each CSV file, preserving the original field names to maintain alignment with the destination schema. The transformation phase included data type standardization, null value handling, and validation against destination constraints.

The loading phase employed two distinct strategies corresponding to the implementation methodologies:

- For Method 1, an incremental approach with automatic schema generation
- For Method 2, a bulk loading process into pre-defined table structures with validation

Error handling protocols were incorporated throughout the pipeline to ensure process resilience, with detailed logging of exceptions encountered during data processing and loading operations.

## **3. Dimensional Modeling for Analytics**

### **3.1 Data Mart Architecture**

A specialized data mart was developed to support analytical processes focused on

LEGO parts. This implementation employed dimensional modeling principles to create an optimized structure for query performance and analytical flexibility.

The dimensional model consists of:

#### **Dimension Tables:**

- Dim\_Color: Contains color attributes including name, RGB value, and transparency
- Dim\_Category: Stores part category classifications

#### **Fact Table:**

- Fact\_Parts: Central fact table containing measures and relationships to dimensions

This star schema design facilitates intuitive data navigation and efficient query execution compared to the normalized operational schema. The model was specifically designed to support analytical requirements related to LEGO part characteristics and their relationships.

The ETL process for the data mart involved:

- Extraction from the normalized operational database
- Transformation to conform data to the dimensional structure
- Loading with appropriate aggregations and surrogate key assignments

This implementation demonstrates effective data warehouse design principles applied to a specialized analytical domain.

## **4. Automated Reporting and Database Maintenance**

### **4.1 Stored Procedure Implementation**

A stored procedure was implemented to generate recurring reports on specific LEGO part characteristics. The procedure design follows database programming best practices:

The procedure encapsulates complex query logic in a centralized, maintainable form, using parameter-based identification of category and color values rather than hard-coded identifiers. This approach enhances maintainability by isolating query implementation from external systems while providing flexibility to accommodate future changes to category or color definitions.

The stored procedure specifically addresses a business requirement to analyze the frequency of Technic Bricks in Dark Bluish Gray color within the inventory. This focused analysis provides valuable insights into product component patterns and distribution across the LEGO product range. The implementation demonstrates several advanced database programming techniques:

First, the procedure utilizes dynamic lookups to retrieve category and color identifiers by name rather than hard-coding these values. This provides resilience against database changes and ensures the procedure continues functioning correctly even if underlying reference data is modified.

Second, the implementation incorporates performance optimization through the use of targeted joins and filtering conditions, minimizing unnecessary data scanning during execution. This approach ensures efficient operation even as the database scales to larger volumes.

Third, result logging capabilities are integrated into the procedure design, creating a historical record of report execution with timestamps and result values in a dedicated ReportLog table. This implementation supports trend analysis and audit requirements while maintaining separation between processing logic and results storage.

The procedure serves as a template for additional analytical routines, establishing a consistent pattern for database-driven reporting that can be extended to other analytical requirements.

## 4.2 Scheduled Automation

The reporting process is automated through SQL Server Agent, with scheduled execution configured for daily runs at 10:00 AM. This implementation leverages database server scheduling capabilities rather than external orchestration tools, creating a self-contained solution with minimal external dependencies.

The SQL Server Agent job configuration employs a modular design pattern that separates concerns across different configuration elements:

- Job definition encapsulates execution parameters and establishes the container for all execution components
- Step definition contains the specific T-SQL command that executes the stored procedure
- Schedule definition controls timing parameters with daily frequency at precisely 10:00 AM
- Server assignment determines the execution environment, binding the job to the local SQL Server instance

This structured approach to automation ensures reliable report generation while providing administrative visibility into the scheduling process. The daily execution timing was specifically chosen to ensure fresh data is available at the beginning of business hours, supporting informed decision-making throughout the workday.

The job implementation includes proper cleanup logic that prevents duplication if deployment scripts are executed multiple times, ensuring environmental consistency. This self-maintaining characteristic reduces administrative overhead and prevents configuration drift over time.

The automation implementation completes the full cycle from data extraction to actionable business insights, creating a closed-loop system that delivers relevant information without manual intervention. The ReportLog table acts as a centralized repository of historical report results, enabling trend analysis and pattern recognition across extended time periods.

## 5. Implementation Advantages and Considerations

### 5.1 Dual-Method Implementation Benefits

The parallel implementation of two distinct methodologies provided valuable insights into database design approaches:

Method 1 (Dynamic Schema) offered advantages in rapid deployment and flexibility, particularly valuable in exploratory data environments where requirements may evolve. It demonstrated effective techniques for post-implementation schema refinement without data loss.

Method 2 (Pre-defined Schema) provided stronger initial data validation and clearer documentation of intended data structures. This approach aligns with traditional data engineering practices that prioritize structural integrity and constraint enforcement.

The comparison between these approaches illustrates important trade-offs in database implementation strategies and provides a foundation for selecting appropriate methodologies based on project requirements.

## **5.2 Performance and Scalability Considerations**

The implementation incorporates several features to address performance and scalability requirements:

The ETL process utilizes chunked data loading patterns to manage memory constraints when processing large datasets. This approach prevents resource exhaustion while maintaining process efficiency.

The dimensional model implementation provides query performance benefits through denormalization and pre-aggregation, reducing join complexity for analytical workloads compared to the normalized operational schema. Indexing strategies focus on foreign key relationships and frequently queried attributes, balancing query performance requirements against storage and maintenance considerations.

## **6. Conclusion**

The LEGO database implementation demonstrates effective application of database design principles, ETL methodologies, and analytical modeling techniques. The dual implementation

approach provides valuable insights into different database development strategies, while the dimensional model and automated reporting capabilities deliver practical analytical value. Future enhancements could include expanded dimensional models, additional automated reports, and integration with business intelligence tools to further leverage the established data architecture.