# Developing a Reusable ETL Framework with Custom Table Format

## Hari Prasad Bomma

Data Engineer, USA
haribomma2007@gmail.com

**Abstract:**

Extracting, transforming, and loading data is a crucial task in data-driven organizations. However, the process of building an ETL workflow can be complex and time-consuming, especially when dealing with diverse data sources and formats. In this research paper we will discuss about a reusable ETL framework that uses custom table format to simplify the development and deployment of ETL processes. This paper will also discuss about the benefits

**Keywords:** Extract, Transform, Load (ETL), Source, Target, Metadata, Workflows, Framework, Data Warehouse

**Introduction:**

A regular ETL (Extract, Transform, Load) framework is essential for managing data flows in businesses. It involves extracting data from various sources, transforming it to fit operational needs, and loading it into a data warehouse or database. The extraction process gathers data from multiple sources such as CRM systems, POS systems, and other databases. During the transformation phase, the data is cleaned, standardized, and enriched to ensure it is accurate and consistent. Finally, the data is loaded into a target system, making it available for analysis and reporting. This process ensures that data is clean, consistent, and ready for analysis.

However, regular ETL frameworks often face several shortages. They can be rigid and difficult to adapt to new data sources or changing business requirements. Metadata changes at the source, such as changes in data types or formats, can be challenging to handle. Additionally, adding new columns or changing data structures requires significant manual intervention. Regular ETL frameworks are also prone to errors when dealing with large volumes of data or complex transformations. They often lack robust error handling and logging mechanisms, making it difficult to track and resolve issues. Furthermore, these frameworks can be resource-intensive, requiring substantial computing power and storage, and may not be easily scalable to accommodate business growth.

Developing a reusable ETL framework is crucial for efficient data management. A reusable ETL framework is modular and flexible, allowing for easy updates and scalability. It automates the ETL process, reducing manual intervention and minimizing errors. With a reusable ETL framework, you can ensure consistent and reliable data processing, which is essential for accurate analysis and reporting.

**Importance of Reusable ETL Frameworks:**

Effective data management is essential for organizations to derive meaningful insights from their data

[2]. ETL processes play a vital role in this, as they are responsible for extracting data from various sources, transforming it into a consistent format, and loading it into a data warehouse or other analytical systems. However, building custom ETL workflow for each data source can be a laborious and error-prone task, especially when dealing with a large number of data sources or frequent changes in the source data structures. [3][1]

To overcome these challenges, the development of reusable ETL frameworks has gained significant attention in the research community. These frameworks aim to provide a standardized and extensible approach to building ETL workflow, reducing the time and effort required for development and maintenance.

**Proposed ETL Framework:**

The ETL framework proposed in this research paper focuses on the use of a custom table format to simplify the development and deployment of ETL processes. The key components of the framework are:

**Custom Table Format:** The framework introduces a custom table format that provides a standardized structure for representing configuration details during the ETL process. By having a standardized table format, the framework ensures consistency and simplifies data handling. It also makes it easier to perform validations, transformations, and error-checking on the data, as the structure is predictable and well-defined.

**Extraction:** The extraction component of the framework is responsible for reading data from various sources, such as relational databases, flat files, or non-relational databases, and converting it into the custom table format. This component can handle diverse data sources and formats, ensuring that data is accurately and efficiently extracted. It also includes mechanisms for incremental data extraction, allowing for periodic updates without the need for a full data load every time.

**Transformation:** The transformation component applies a series of data manipulation operations, such as filtering, merging, or enriching the data, to ensure that it meets the required quality standards and aligns with the target data model. These operations can include data cleaning, standardization, aggregation, and business rule application. By transforming the data, the framework ensures that it is consistent, accurate, and ready for analysis or reporting, providing meaningful insights to stakeholders.

**Loading:** The loading component is responsible for inserting the transformed data into the target data warehouse or analytical system. This step involves writing the cleaned and processed data into the designated storage location. The loading process is optimized for performance, ensuring that large volumes of data can be loaded quickly and efficiently. It also includes mechanisms for handling data overwrites, updates, and deletes, maintaining data integrity and consistency.

**Metadata Management:** The framework includes a metadata management component that stores information about the ETL processes, data sources, and transformations. This metadata is crucial for tracking data lineage, auditing, and compliance. It provides a clear understanding of how data flows through the ETL process, enabling easier troubleshooting and maintenance. Additionally, metadata management helps in seamless updates and modifications to the ETL processes.

**Reusability and Extensibility:** The framework is designed to be highly reusable and extensible; allowing developers to easily integrate new data sources or modify existing ETL processes without extensive rework. This modular design ensures that components can be reused across different projects and jobs, reducing development time and effort. Extensibility allows the framework to adapt to changing business requirements and data sources, ensuring it remains relevant and effective over time.

**Benefits and Limitations:**

**Benefits:**

The use of a custom table format in the proposed ETL framework offers several benefits as well as its own limitations.

**Standardized Data Representation:** The custom table format provides a consistent way to represent configurations and queries during the ETL process, making it easier to integrate and process data from diverse sources [4].

**Simplified Development:** The standardized format simplifies the development of ETL processes, as developers can focus on the data transformation logic rather than the complexities of different data source formats.

**Improved Maintainability:** The reusable nature of the framework and the use of a custom table format make it easier to maintain and update ETL processes, reducing the time and effort required for future enhancements or modifications.

**Enhanced Scalability:** The framework's modular design and the use of a custom table format allow for the easy integration of new data sources or the expansion of existing ETL processes, enabling organizations to scale their data management capabilities as their needs evolve.

**Reduced Errors:** By providing a standardized approach to ETL, the proposed framework can help reduce the risk of errors and inconsistencies that may arise from manual data processing or the use of multiple disparate ETL tools.

**Increased Efficiency:** The streamlined development and deployment process enabled by the custom table format and the reusable framework can lead to significant improvements in the overall efficiency of the ETL process.

**Limitations:**

**Complexity**: The initial setup and configuration of the custom table format and metadata management can be complex. Requires careful planning and design to ensure all components work seamlessly together.

**Dependency Management**: Ensuring all dependencies are properly managed and maintained can be challenging. Changes in one component may have cascading effects on other parts of the framework.

**Maintenance and Updates**: Regular maintenance and updates are required to keep the framework functioning optimally. Requires dedicated resources to monitor, troubleshoot, and enhance the framework as needed.

**Case Study:** Developing a Reusable ETL Framework with Custom Table Format

**Objective**

The goal is to create a reusable ETL framework that can handle frequent metadata changes by using custom tables to manage source and target connections, queries, table names, schemas, and mappings.

**Custom Table Format**

We start by defining a custom table format to provide a standardized structure for representing data during the ETL process. This includes:

- **Source Connection**: Information about the source database connection, such as host, port, username, password, and database name.
- **Source Query**: The SQL query used to extract data from the source.
- **Source Table Name**: The name of the table in the source database.

- **Source Schema**: The schema of the source table.
- **Target Connection**: Information about the target database connection.
- **Target Query**: The SQL query used to insert data into the target.
- **Target Table Name**: The name of the table in the target database.
- **Target Schema**: The schema of the target table.
- **Load Type**: Incremental /Full
- **Incremental Lookup:** column to find incremental delta in source system.
- **Last Processed Status:** 0/1 to represent success or failure
- **Active Indicator:** To maintain or deactivate the record in custom table
- **Last Processed:** To capture when was the data load last processed
- **Last Updated:** To capture when this record got last updated
- **Archival required:** 0/1 to represent archival required /not after the data load complete

## Data Extraction

The extraction component reads data from various sources, such as relational databases, flat files, or non-relational databases, and converts it into the custom table format. The framework uses the source connection details to connect to the source database and execute the source query to extract the required data.

**Data Loading:** The loading component inserts the transformed data into the target data warehouse or analytical system. The framework uses the target connection details to connect to the target database and execute the target query to load the data. This step ensures that the data is accurately and efficiently loaded into the target system.

**Metadata Management:** The framework includes a metadata management component that stores information about the ETL processes, data sources, and transformations. This component maintains details about source and target connections, queries, table names, and schemas. It also tracks changes in metadata, allowing the framework to adapt to new data sources or modifications in existing sources.

**Custom Table for Mappings:** In addition to the custom table format for source and target details, we introduce another custom table for mappings from source to target. This table includes:

- **Source Column Name**: The name of the column in the source table.
- **Target Column Name**: The name of the column in the target table.
- **Transformation Rules**: Any transformation rules that need to be applied to the source column before loading it into the target column.

## Results:

**Dynamic Value Passing:** Values can be dynamically passed from the custom table, enabling seamless execution of the ETL process. This approach reduces the need for hardcoding values in the ETL scripts and allows for greater flexibility and adaptability. By dynamically fetching values from the custom table, the framework can easily adjust to changes in source or target configurations without modifying the core ETL logic.

**Centralized Changes:** Changes can be managed in one table rather than updating all affected tables and dependencies, streamlining the update process. This centralization minimizes the risk of errors and inconsistencies across different parts of the ETL pipeline. By maintaining a single source of truth for configurations and mappings, updates can be implemented more efficiently and with fewer disruptions to the overall workflow.

**Handling Metadata Changes:** Supports frequent metadata changes, such as the addition of new columns or removal of existing columns. This capability allows the ETL framework to quickly adapt to evolving data requirements without extensive rework. When new columns are added to the source data, the corresponding updates can be made directly in the custom table, ensuring that the ETL process continues to function smoothly.

**Modular and Flexible Design:** The reusable ETL framework is designed to be modular and flexible, ensuring easy updates and scalability. Each component of the ETL process (extraction, transformation, loading, and metadata management) can be independently developed, tested, and maintained. This modularity not only simplifies the development process but also makes it easier to scale the framework to accommodate

growing data volumes and new data sources. By leveraging a flexible design, the framework can be extended with new features and functionalities as needed, ensuring it remains effective over time.

**Conclusion:**

In conclusion, developing a reusable ETL framework with a custom table format significantly enhances data management efficiency. By dynamically passing values from the custom table, the ETL process becomes seamless and adaptable to changes. Centralizing changes in one table simplifies updates and reduces the risk of errors and inconsistencies. The framework's ability to handle frequent metadata changes, such as the addition of new columns, ensures that the ETL process remains flexible and can quickly adapt to evolving data requirements. Moreover, the modular and flexible design of the framework allows for easy updates and scalability, making it capable of handling growing data volumes and new data sources. This approach ensures reliable and consistent data processing, providing a robust foundation for accurate analysis and reporting. Ultimately, a reusable ETL framework not only streamlines data management but also supports the long-term goals of scalability and adaptability in a dynamic business environment.

**References:**

1. S. Vyas and P. Vaishnav, "A comparative study of various ETL process and their testing techniques in data warehouse," Jul. 04, 2017, Taylor & Francis. doi: 10.1080/09720510.2017.1395194.
2. P. Ponniah, "Data Extraction, Transformation, and Loading." p. 257, Aug. 24, 2001. doi: 10.1002/0471221627.ch12.
3. S. El–Sappagh, A. Hendawi, and A. H. E. Bastawissy, "A proposed model for data warehouse ETL processes," May 11, 2011, Elsevier BV. doi: 10.1016/j.jksuci.2011.05.005.
4. R. Venkatakrishnan, "Design, Implementation, and Assessment of Innovative Data Warehousing; Extract, Transformation, and Load(ETL); and Online Analytical Processing(OLAP) on BI," Jun. 30, 2020. doi: 10.5121/ijdms.2020.12301.
5. A. Schwinn and J. Schelp, "Data Integration Patterns," 2003. [Online]. Available: https://www.alexandria.unisg.ch/213454/