

Extract - Transform - Load

Time Series Team

Ngày 12 tháng 8 năm 2025

Mục lục

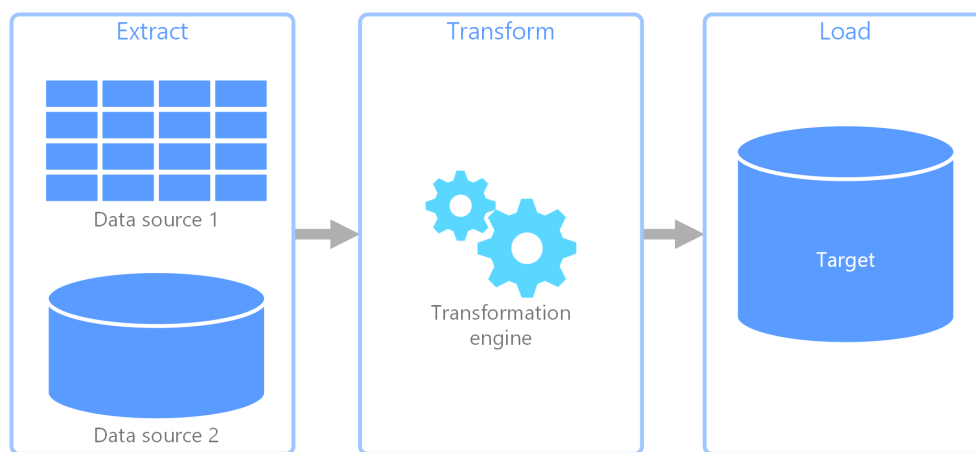
1	1. Tổng quan về ETL Pipeline	2
1.1	ETL là gì?	2
1.2	Tại sao cần ETL Pipeline?	2
1.3	Kiến trúc dữ liệu hiện đại	2
1.4	Kiến trúc dữ liệu hiện đại	3
1.4.1	Mô hình Business Intelligence Pyramid	3
1.4.2	Quy trình ETL tổng quát	4
1.4.3	Data Warehouse như một nhà hàng	5
1.4.4	Lợi ích của kiến trúc ETL	6
2	2. Data Extraction - Trích xuất dữ liệu	7
2.1	Quy trình Extraction	7
2.2	Case Study: Iowa Liquor Sales Dataset	8
2.3	Xử lý vấn đề Composite Key Collisions	8
2.4	Batch Processing	9
3	3. Data Transformation - Chuyển đổi dữ liệu	9
3.1	Data Cleaning - Làm sạch dữ liệu	10
3.1.1	Xử lý dữ liệu trùng lặp	10
3.1.2	Xử lý dữ liệu thiếu (Missing Data)	11
3.1.3	Xử lý dữ liệu không hợp lệ	11
3.2	Type Casting - Chuyển đổi kiểu dữ liệu	11
3.3	Data Processing - Xử lý dữ liệu	12
3.3.1	Derived Columns - Tạo cột phái sinh	12
4	4. Data Loading - Tải dữ liệu	12
4.1	Dimension và Fact Tables	12
4.1.1	Thiết kế Star Schema	12
4.2	Slowly Changing Dimensions (SCD)	13
4.3	Tại sao cần quan tâm đến Slowly Changing Dimensions?	14
4.4	Các loại Slowly Changing Dimensions phổ biến	14
4.5	Ví dụ minh họa SCD Type 2	14
4.6	Ví dụ minh họa các SCD khác	15
4.6.1	SCD Type 1 - Overwrite	15
4.6.2	SCD Type 2 - Add New Row	15
5	5. Advanced ETL Concepts	16
5.1	Change Data Capture (CDC)	16
5.2	Data Quality Monitoring	16
5.3	Performance Optimization	17
5.3.1	Parallel Processing	17
5.3.2	Memory Optimization	18
5.4	Error Handling và Monitoring	18

1. Tổng quan về ETL Pipeline

1.1 ETL là gì?

ETL (Extract, Transform, Load) là quy trình quan trọng trong Data Engineering, bao gồm ba giai đoạn chính:

- **Extract (Trích xuất):** Thu thập dữ liệu từ nhiều nguồn khác nhau
- **Transform (Chuyển đổi):** Làm sạch, xử lý và chuẩn hóa dữ liệu
- **Load (Tải):** Lưu trữ dữ liệu đã xử lý vào Data Warehouse



Hình 1: Extract-Transform-Load

1.2 Tại sao cần ETL Pipeline?

Trong môi trường doanh nghiệp hiện đại, việc ra quyết định dựa trên dữ liệu (Data-driven Decision Making) đã trở thành yếu tố then chốt. Tuy nhiên, các tổ chức thường gặp phải những thách thức sau:

1. **Tốn thời gian:** Xây dựng báo cáo từ đầu, thu thập dữ liệu từ nhiều nguồn
2. **Nhiều người tham gia:** Dẫn đến sự không nhất quán và lỗi con người
3. **Khó tích hợp báo cáo:** Nhiều công cụ, nhiều thuật ngữ khác nhau
4. **Vấn đề không lường trước:** Số liệu không khớp, thông tin sai lệch

1.3 Kiến trúc dữ liệu hiện đại

ETL Pipeline là nền tảng của Business Intelligence, nằm trong hệ thống phân cấp:

- **Data Sourcing:** Thu thập dữ liệu từ các nguồn
- **Data Warehousing (ETL):** Xử lý và lưu trữ dữ liệu
- **Data Mining:** Khai thác thông tin từ dữ liệu
- **Report Data Analysis:** Phân tích và báo cáo
- **Decision Making:** Ra quyết định dựa trên dữ liệu

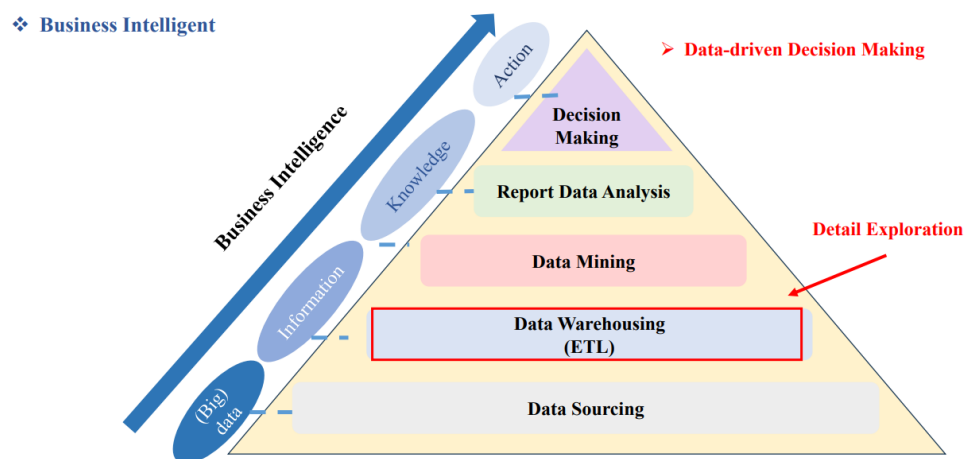
1.4 Kiến trúc dữ liệu hiện đại

1.4.1 Mô hình Business Intelligence Pyramid

Business Intelligence (BI) được tổ chức theo mô hình kim tự tháp, trong đó ETL Pipeline đóng vai trò nền tảng quan trọng. Mũi tên "Business Intelligence" biểu thị sự gia tăng giá trị và độ phức tạp từ dữ liệu thô đến knowledge và cuối cùng là action. Càng lên cao trong pyramid, chúng ta có:

- **Tăng giá trị business:** Từ raw data đến actionable insights
- **Giảm volume:** Từ massive datasets đến focused recommendations
- **Tăng context:** Từ isolated data points đến business-relevant information
- **Tăng impact:** Từ descriptive reports đến predictive và prescriptive analytics

Mô hình này minh họa sự tiến hóa từ dữ liệu thô đến insight và hành động cụ thể:



Hình 2: Business Intelligence Pyramid

Các tầng trong BI Pyramid:

1. Big Data (Tầng nền tảng):

- Dữ liệu thô từ nhiều nguồn khác nhau
- Volume lớn, đa dạng về format và structure
- Yêu cầu khả năng lưu trữ và xử lý mạnh mẽ

2. Data Sourcing (Thu thập dữ liệu):

- Identification và connection đến các data sources
- Data profiling và quality assessment
- Establish data access patterns và security

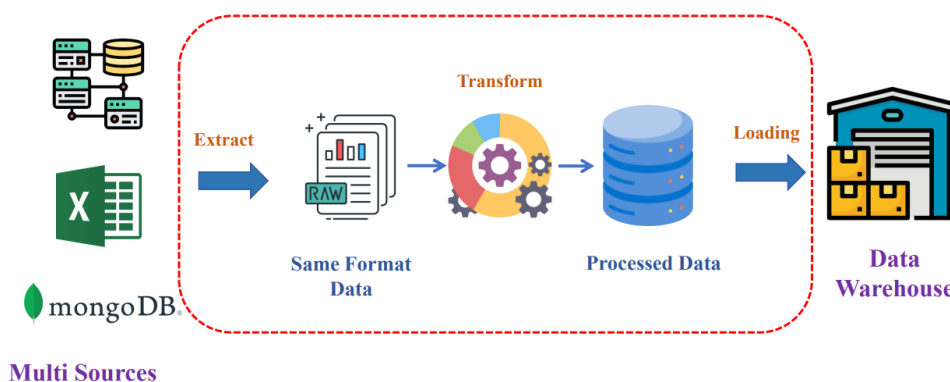
3. Data Warehousing - ETL (Xử lý và lưu trữ):

- Extract, Transform, Load processes
- Data integration và standardization

- Historical data storage và management
4. **Data Mining (Khai thác dữ liệu):**
 - Pattern recognition và anomaly detection
 - Statistical analysis và machine learning
 - Predictive modeling và forecasting
 5. **Report Data Analysis (Phân tích và báo cáo):**
 - Descriptive analytics và visualization
 - KPI monitoring và dashboard creation
 - Ad-hoc analysis và self-service BI
 6. **Decision Making (Ra quyết định):**
 - Data-driven decision processes
 - Strategic planning dựa trên insights
 - Action planning và execution

1.4.2 Quy trình ETL tổng quát

ETL Pipeline là connecting bridge giữa raw data sources và business intelligence applications:



Hình 3: ETL Process Flow

Multi Sources (Đa nguồn dữ liệu):

- **Relational Databases:** MySQL, PostgreSQL, SQL Server, Oracle
- **File Systems:** Excel, CSV, JSON, XML, Parquet
- **NoSQL Databases:** MongoDB, Cassandra, Redis
- **Cloud Storage:** AWS S3, Google Cloud Storage, Azure Blob
- **APIs:** REST APIs, GraphQL, SOAP services
- **Streaming Data:** Kafka, Kinesis, Event Hubs

Extract Phase (Trích xuất):

- Kết nối và đọc dữ liệu từ multiple sources
- Handle different data formats và protocols
- Implement extraction strategies (full, incremental, CDC)
- Data validation và error handling

Transform Phase (Chuyển đổi):

- **Same Format Data:** Standardization về common format
- Data cleaning và quality improvement
- Business rule application
- Schema mapping và data type conversion
- Aggregation và calculation

Load Phase (Tải dữ liệu):

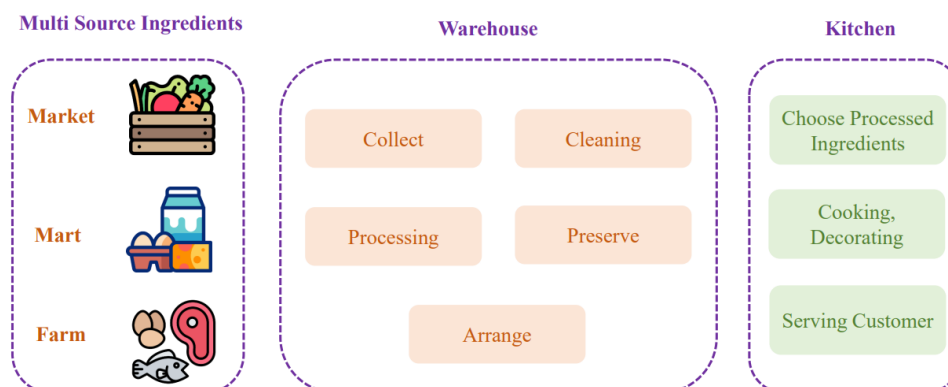
- **Processed Data:** Optimized cho analytical workloads
- Dimensional modeling (Star/Snowflake schema)
- Indexing và partitioning cho performance
- Data archival và retention policies

Data Warehouse (Kho dữ liệu):

- Central repository cho enterprise data
- Optimized cho OLAP operations
- Support cho reporting và analytics tools
- Data governance và security implementation

1.4.3 Data Warehouse như một nhà hàng

Để hiểu rõ hơn về vai trò của Data Warehouse, chúng ta có thể tưởng tượng nó như một nhà hàng chuyên nghiệp:



Hình 4: Data Warehouse as Restaurant

Multi-Source Ingredients (Nguyên liệu đa nguồn):

Giống như một nhà hàng cần nguyên liệu từ nhiều nguồn khác nhau:

- **Excel:** Như market stalls - dữ liệu nhỏ lẻ, manual
- **MySQL/PostgreSQL:** Như wholesale suppliers - dữ liệu structured, reliable
- **Supabase:** Như modern food distributors - cloud-based, scalable
- **APIs:** Như direct farm connections - real-time, fresh data

Data Warehouse (Kho lưu trữ):

- **Collect:** Thu thập nguyên liệu từ các nguồn
- **Cleaning:** Làm sạch, loại bỏ impurities
- **Processing:** Chế biến theo standardized recipes (business rules)
- **Preserve:** Bảo quản trong điều kiện optimal
- **Arrange:** Sắp xếp theo organized structure

Serving Customers (Phục vụ người dùng):

Data Warehouse phục vụ nhiều types of "customers":

- **Business Analysts:** Như regular diners - cần standard reports
- **Data Scientists:** Như professional chefs - cần raw ingredients cho custom analysis
- **Executives:** Như VIP customers - cần high-level dashboards và insights
- **Operational Teams:** Như catering orders - cần scheduled, automated deliveries

Quality Control (Kiểm soát chất lượng):

- **Data Quality:** Như food safety standards
- **Performance Monitoring:** Như service time tracking
- **Customer Satisfaction:** Như feedback collection từ users
- **Continuous Improvement:** Như menu updates và process optimization

1.4.4 Lợi ích của kiến trúc ETL

Business Benefits:

1. **Single Source of Truth:** Centralized data repository
2. **Data Consistency:** Standardized definitions và calculations
3. **Historical Analysis:** Time-series data cho trend analysis
4. **Improved Performance:** Optimized cho analytical queries
5. **Data Governance:** Centralized security và access control

Technical Benefits:

1. **Scalability:** Handle growing data volumes
2. **Reliability:** Robust error handling và recovery
3. **Maintainability:** Modular design với clear separation of concerns
4. **Monitoring:** Comprehensive logging và alerting
5. **Flexibility:** Support cho multiple output formats và destinations

Với foundation này, ETL Pipeline trở thành backbone của modern data-driven organization, enabling chuyển đổi từ reactive reporting sang proactive, insight-driven decision making.

2. Data Extraction - Trích xuất dữ liệu

2.1 Quy trình Extraction

Data Extraction (Trích xuất dữ liệu) là giai đoạn đầu tiên trong quy trình ETL (Extract - Transform - Load), đảm nhiệm việc lấy dữ liệu từ các nguồn khác nhau như cơ sở dữ liệu, API, file logs, hoặc các hệ thống ứng dụng. Quá trình này cần đảm bảo tính đầy đủ, chính xác và kịp thời của dữ liệu để phục vụ cho các bước xử lý tiếp theo. Các phương pháp trích xuất dữ liệu phổ biến gồm:

- **Pull (Active):** Hệ thống ETL chủ động gửi yêu cầu đến nguồn dữ liệu theo lịch định sẵn hoặc theo sự kiện để *kéo* dữ liệu về.
 - *Ưu điểm:* Kiểm soát được thời điểm và tần suất lấy dữ liệu; dễ tích hợp với nhiều loại nguồn.
 - *Nhược điểm:* Có thể gây tải cao lên hệ thống nguồn nếu tần suất truy vấn lớn; độ trễ phụ thuộc vào lịch lấy dữ liệu.
 - *Ví dụ:* Hệ thống ETL chạy cron job mỗi 1 giờ để truy vấn bảng giao dịch từ cơ sở dữ liệu bán hàng.
- **Push (Inactive):** Nguồn dữ liệu chủ động *đẩy* dữ liệu đến hệ thống ETL khi có dữ liệu mới hoặc khi đạt điều kiện nhất định.
 - *Ưu điểm:* Độ trễ thấp; giảm tải cho hệ thống ETL vì không cần truy vấn liên tục.
 - *Nhược điểm:* Cần cấu hình và hỗ trợ từ hệ thống nguồn; khó kiểm soát nếu nguồn gửi dữ liệu sai định dạng.
 - *Ví dụ:* Ứng dụng thương mại điện tử gửi sự kiện đơn hàng mới qua Webhook đến hệ thống ETL ngay khi đơn được tạo.
- **Incremental:** Chỉ trích xuất phần dữ liệu mới hoặc thay đổi so với lần trích xuất trước đó, thường dựa trên các cột như `last_modified` hoặc `created_at`.
 - *Ưu điểm:* Tiết kiệm thời gian và tài nguyên; giảm dung lượng dữ liệu cần xử lý.
 - *Nhược điểm:* Cần cơ chế theo dõi trạng thái lần trích xuất cuối; dễ bỏ sót nếu cột thời gian không được cập nhật chính xác.
 - *Ví dụ:* ETL chỉ lấy các bản ghi giao dịch có `updated_at` lớn hơn thời điểm trích xuất gần nhất.
- **Change Data Capture (CDC):** Theo dõi và ghi lại mọi thay đổi (INSERT, UPDATE, DELETE) trên dữ liệu nguồn bằng log hoặc trigger ở mức cơ sở dữ liệu.

- *Ưu điểm*: Ghi nhận chính xác và đầy đủ các thay đổi; phù hợp cho đồng bộ dữ liệu thời gian thực.
- *Nhược điểm*: Cấu hình phức tạp; có thể ảnh hưởng hiệu năng của hệ thống nguồn nếu không tối ưu.
- *Ví dụ*: Sử dụng Debezium để lắng nghe binlog của MySQL, phát hiện mọi thay đổi và gửi đến Kafka.

2.2 Case Study: Iowa Liquor Sales Dataset

Chúng ta sẽ sử dụng bộ dữ liệu Iowa Liquor Sales làm ví dụ thực tế. Dataset này chứa thông tin về việc mua bán rượu tại Iowa với các trường chính:

```

1 # Data Structure of Iowa Liquor Sales
2 columns = [
3     'invoice_line_no',
4     'date',
5     'store',
6     'name',
7     'address',
8     'city',
9     'zipcode',
10    'store_location',
11    'county_number',
12    'county',
13    'category',
14    'category_name',
15    'vendor_no',
16    'vendor_name',
17    'item_number',
18    'item_description',
19    'pack',
20    'bottle_volume_ml',
21    'state_bottle_cost',
22    'state_bottle_retail',
23    'sale_bottles',
24    'sale_dollars',
25    'sale_liters',
26    'sale_gallons'
27 ]

```

2.3 Xử lý vấn đề Composite Key Collisions

Khi extract dữ liệu từ nhiều nguồn, thường gặp vấn đề trùng lặp composite key:

```

1 -- Create staging table to extract data
2 CREATE TABLE Staging_Sales (
3     invoice_line_no TEXT,
4     date TEXT,
5     store TEXT,
6     name TEXT,
7     address TEXT,
8     city TEXT,
9     zipcode TEXT,
10    store_location TEXT,
11    county_number TEXT,
12    county TEXT,
13    category TEXT,

```

```

14     category_name TEXT,
15     vendor_no TEXT,
16     vendor_name TEXT,
17     item_number TEXT,
18     item_description TEXT,
19     pack TEXT,
20     bottle_volume_ml TEXT,
21     state_bottle_cost TEXT,
22     state_bottle_retail TEXT,
23     sale_bottles TEXT,
24     sale_dollars TEXT,
25     sale_liters TEXT,
26     sale_gallons TEXT,
27     file_name TEXT NOT NULL,
28     processed_timestamp DATETIME NOT NULL,
29     record_source TEXT NOT NULL
30 );

```

2.4 Batch Processing

Để xử lý hiệu quả, dữ liệu được chia thành các batch nhỏ:

```

1 def process_csv_in_batches(file_path, batch_size=10000):
2     """
3     Process CSV by batch_size
4     """
5     import pandas as pd
6
7     chunk_iter = pd.read_csv(file_path, chunksize=batch_size)
8
9     for i, chunk in enumerate(chunk_iter):
10        # Add metadata
11        chunk['file_name'] = file_path
12        chunk['processed_timestamp'] = pd.Timestamp.now()
13        chunk['record_source'] = 'batch_processing'
14
15        # Save into staging area
16        chunk.to_sql('staging_sales', conn,
17                    if_exists='append', index=False)
18
19        print(f"Processed batch {i+1}: {len(chunk)} records")

```

3. Data Transformation - Chuyển đổi dữ liệu

Chuyển đổi dữ liệu (Data Transformation) là giai đoạn quan trọng trong quy trình xử lý dữ liệu, nhằm làm sạch và biến đổi dữ liệu thô thành dữ liệu có chất lượng và phù hợp cho các bước phân tích hoặc lưu trữ. Quá trình này thường bao gồm hai bước chính: **Data Cleaning (Làm sạch dữ liệu)** và **Data Processing (Xử lý dữ liệu)**.

Data Cleaning - Làm sạch dữ liệu

Trong bước này, dữ liệu được kiểm tra và xử lý để loại bỏ hoặc sửa các lỗi và vấn đề ảnh hưởng đến chất lượng dữ liệu như:

- **Duplicated (Trùng lặp):** Xác định và loại bỏ các bản ghi dữ liệu bị trùng để tránh gây sai lệch trong phân tích.

- **Filtering (Lọc dữ liệu):** Lọc bỏ các dữ liệu không hợp lệ hoặc không cần thiết dựa trên điều kiện cụ thể.
- **Missing (Dữ liệu thiếu):** Xử lý các giá trị bị thiếu bằng cách điền giá trị thay thế, loại bỏ hoặc giữ nguyên tùy theo mục đích.
- **Invalid (Dữ liệu không hợp lệ):** Phát hiện và sửa các dữ liệu sai định dạng, nằm ngoài phạm vi cho phép hoặc không hợp lệ.
- **Type Casting (Chuyển đổi kiểu dữ liệu):** Đảm bảo các trường dữ liệu có kiểu dữ liệu phù hợp (ví dụ: số nguyên, chuỗi, ngày tháng) để dễ dàng xử lý và phân tích.

Data Processing - Xử lý dữ liệu

Sau khi dữ liệu đã được làm sạch, bước xử lý dữ liệu sẽ biến đổi dữ liệu để phục vụ cho các mục tiêu phân tích, bao gồm:

- **Enrichment (Làm giàu dữ liệu):** Bổ sung thêm thông tin hoặc thuộc tính mới cho dữ liệu hiện có để tăng giá trị khai thác.
- **Aggregate (Tổng hợp):** Tập hợp dữ liệu theo các nhóm hoặc tiêu chí nhất định để tạo ra các chỉ số tổng quan, ví dụ như tổng doanh số theo tháng.
- **Normalize (Chuẩn hóa):** Điều chỉnh dữ liệu về cùng một thang đo hoặc định dạng để dễ so sánh và phân tích.
- **Integrate (Tích hợp):** Kết hợp dữ liệu từ nhiều nguồn khác nhau thành một bộ dữ liệu thống nhất.
- **Derived Column (Cột dẫn xuất):** Tạo thêm các cột dữ liệu mới dựa trên phép tính hoặc kết hợp các cột hiện có, ví dụ như tính tổng giá trị hoặc phân loại nhóm.

Nhờ vào các bước làm sạch và xử lý này, dữ liệu được chuyển đổi trở nên đồng nhất, chính xác và có ý nghĩa hơn, từ đó giúp nâng cao hiệu quả trong việc phân tích và ra quyết định.

3.1 Data Cleaning - Làm sạch dữ liệu

3.1.1 Xử lý dữ liệu trùng lặp

```
1 import pandas as pd
2
3 # Find and Solve duplicate
4 df_full_duplicated = df[df.duplicated(keep=False)]
5 print(f"No. duplicate records: {len(df_full_duplicated)}")
6
7 # Drop duplicate
8 df_cleaned = df.drop_duplicates()
9
10 # Solve duplicate by key
11 df_special_duplicated = df[df.duplicated(
12     subset=['invoice_line_no', 'store'], keep=False
13 )]
14
15 df_final = df[
16     df.duplicated(subset=['invoice_line_no', 'store'], keep='first') |
17     ~df.duplicated(subset=['invoice_line_no', 'store'], keep=False)
18 ]
```

3.1.2 Xử lý dữ liệu thiếu (Missing Data)

```
1 # Check null data
2 null_columns = df.columns[df.isnull().sum() > 0].tolist()
3 print("Columns having null value:", null_columns)
4
5 # Fill missing data
6 strategies = {
7     'address': 'Unknown',
8     'city': 'Unknown',
9     'zipcode': 'Unknown',
10    'county_number': 'Unknown',
11    'county': 'Unknown',
12    'category': 'Unknown',
13    'category_name': 'Unknown'
14 }
15
16 df_filled = df.fillna(value=strategies, inplace=True)
17
18 # Delete columns that have missing values
19 critical_columns = [
20     'state_bottle_cost', 'state_bottle_retail',
21     'sale_bottles', 'sale_dollars', 'sale_liters', 'sale_gallons'
22 ]
23
24 # Delete columns based on business logic
25 df_cleaned = df.dropna(subset=critical_columns)
```

3.1.3 Xử lý dữ liệu không hợp lệ

```
1 # Check numeric columns
2 numeric_cols = [
3     'pack', 'bottle_volume_ml', 'state_bottle_cost',
4     'state_bottle_retail', 'sale_bottles', 'sale_dollars',
5     'sale_liters', 'sale_gallons'
6 ]
7 for col in numeric_cols:
8     df[col] = pd.to_numeric(df[col], errors='coerce')
9
10 # Delete invalid data (< 1)
11 df_valid = df[(df[numeric_cols] >= 1).all(axis=1)]
```

3.2 Type Casting - Chuyển đổi kiểu dữ liệu

```
1 numeric_cols = [
2     'state_bottle_cost', 'state_bottle_retail', 'sale_bottles',
3     'sale_dollars', 'sale_liters', 'sale_gallons'
4 ]
5
6 for col in numeric_cols:
7     df[col] = pd.to_numeric(df[col], errors='coerce')
8
9 # Change data type into datetime
10 df['date'] = pd.to_datetime(df['date'], errors='coerce')
```

3.3 Data Processing - Xử lý dữ liệu

3.3.1 Derived Columns - Tạo cột phái sinh

```

1 # Create columns based on business logic
2 df['revenue'] = df['sale_dollars']
3 df['cost'] = df['state_bottle_cost'] * df['sale_bottles']
4 df['profit'] = (df['state_bottle_retail'] - df['state_bottle_cost']) * df['
    sale_bottles']
5 df['total_bottles_sold'] = df['sale_bottles']
6 df['total_volume_sold_in_liters'] = df['sale_liters']
7
8 # Calculate advance metrics
9 df['profit_margin'] = (
10     (df['profit'] / df['revenue'] * 100)
11     .round(2)
12     .where(df['revenue'] > 0, 0)
13 )
14
15 df['average_bottle_price'] = (
16     (df['sale_dollars'] / df['sale_bottles'])
17     .round(2)
18     .where(df['sale_bottles'] > 0, 0)
19 )
20
21 df['volume_per_bottle_sold'] = (
22     (df['sale_liters'] / df['sale_bottles'])
23     .round(2)
24     .where(df['sale_bottles'] > 0, 0)
25 )

```

4. Data Loading - Tải dữ liệu

4.1 Dimension và Fact Tables

4.1.1 Thiết kế Star Schema

Trong Data Warehouse, chúng ta sử dụng mô hình Star Schema với:

Fact Table (sales_fact):

```

1 CREATE TABLE sales_fact (
2     sales_key INTEGER PRIMARY KEY,
3     invoice_line_no TEXT,
4     store INTEGER,
5     date_key INTEGER,
6     store_key INTEGER,
7     item_key INTEGER,
8     vendor_key INTEGER,
9     revenue DECIMAL,
10    profit DECIMAL,
11    cost DECIMAL,
12    total_bottles_sold INTEGER,
13    total_volume_sold_in_liters DECIMAL,
14    profit_margin DECIMAL,
15    average_bottle_price DECIMAL,
16    volume_per_bottle_sold DECIMAL,
17    processed_timestamp DATETIME
18 );

```

Dimension Tables:

```

1  -- Date Dimension
2  CREATE TABLE date_dim (
3      date_key INTEGER PRIMARY KEY,
4      date DATE,
5      year INTEGER,
6      month INTEGER,
7      day INTEGER,
8      quarter INTEGER,
9      weekday TEXT
10 );
11
12 -- Store Dimension
13 CREATE TABLE store_dim (
14     store_key INTEGER PRIMARY KEY,
15     store_id INTEGER,
16     address TEXT,
17     city TEXT,
18     zipcode TEXT,
19     store_location TEXT,
20     county_number TEXT,
21     county TEXT,
22     start_date DATE,
23     end_date DATE,
24     is_active BOOLEAN
25 );
26
27 -- Item Dimension
28 CREATE TABLE item_dim (
29     item_key INTEGER PRIMARY KEY,
30     item_number TEXT,
31     item_description TEXT,
32     category TEXT,
33     category_name TEXT,
34     pack FLOAT,
35     bottle_volume_ml FLOAT,
36     state_bottle_cost DECIMAL,
37     state_bottle_retail DECIMAL,
38     start_date DATE,
39     end_date DATE,
40     is_active BOOLEAN
41 );
42
43 -- Vendor Dimension
44 CREATE TABLE vendor_dim (
45     vendor_key INTEGER PRIMARY KEY,
46     vendor_no TEXT,
47     vendor_name TEXT,
48     start_date DATE,
49     end_date DATE,
50     is_active BOOLEAN
51 );

```

4.2 Slowly Changing Dimensions (SCD)

Slowly Changing Dimensions (SCD) là khái niệm dùng để mô tả các chiều dữ liệu trong kho dữ liệu mà thông tin trong các chiều này thay đổi không thường xuyên hoặc thay đổi chậm theo thời gian. Ví dụ về các chiều như vậy thường là thông tin khách hàng, sản phẩm, địa điểm, nhân viên, v.v.

Trong quá trình quản lý dữ liệu, việc xử lý các chiều này rất quan trọng vì chúng không chỉ đơn thuần chứa dữ liệu hiện tại mà còn cần giữ được lịch sử các thay đổi theo thời gian để phục vụ phân tích và báo cáo chính xác.

4.3 Tại sao cần quan tâm đến Slowly Changing Dimensions?

- Các chiều dữ liệu thường lưu các thuộc tính mô tả đối tượng như tên, địa chỉ, trạng thái, hoặc các thuộc tính khác.
- Khi các thuộc tính này thay đổi (ví dụ khách hàng đổi địa chỉ, nhân viên thay đổi chức vụ), ta cần quyết định:
 - Có giữ lại thông tin cũ không?
 - Nếu giữ thì bằng cách nào?
 - Nếu thay đổi thì có ảnh hưởng đến dữ liệu báo cáo không?

4.4 Các loại Slowly Changing Dimensions phổ biến

1. SCD Type 0 (Không thay đổi)

Không bao giờ cập nhật dữ liệu cũ, giữ nguyên giá trị ban đầu. Thường dùng khi không cần theo dõi lịch sử thay đổi.

2. SCD Type 1 (Ghi đè dữ liệu cũ)

Khi có thay đổi, dữ liệu cũ bị ghi đè bởi dữ liệu mới, không giữ lịch sử thay đổi. Ví dụ: sửa lỗi chính tả trong tên khách hàng.

3. SCD Type 2 (Giữ lịch sử đầy đủ)

Tạo bản ghi mới mỗi khi có thay đổi, giữ lại lịch sử của các bản ghi cũ. Thường dùng các trường như `start_date`, `end_date` hoặc `is_current` để xác định phiên bản hiện tại. Đây là loại SCD phổ biến nhất.

4. SCD Type 3 (Lưu một số phiên bản thay đổi)

Chỉ lưu một số thuộc tính cũ và mới trong cùng một bản ghi, không giữ lịch sử đầy đủ. Ví dụ, có thể lưu giá trị cũ và giá trị hiện tại của một thuộc tính cụ thể.

5. Các loại SCD mở rộng khác (Type 4, Type 6, ...)

Kết hợp các phương pháp trên hoặc lưu lịch sử trong bảng riêng biệt.

4.5 Ví dụ minh họa SCD Type 2

Giả sử bảng khách hàng lưu thông tin địa chỉ. Khi khách hàng thay đổi địa chỉ, thay vì cập nhật trực tiếp trên bản ghi hiện tại, ta tạo thêm một bản ghi mới với địa chỉ mới và đánh dấu bản ghi cũ là không còn hiệu lực (ví dụ `end_date` = ngày thay đổi).

CustomerID	Name	Address	Start_Date	End_Date	Is_Current
1001	Nguyễn An	123 Đường A	2020-01-01	2022-06-30	False
1001	Nguyễn An	456 Đường B	2022-07-01	NULL	True

Bảng 1: Ví dụ SCD Type 2 cho quản lý địa chỉ khách hàng

Điều này giúp khi phân tích doanh số theo thời gian, ta biết chính xác khách hàng ở đâu trong từng giai đoạn.

4.6 Ví dụ minh họa các SCD khác

4.6.1 SCD Type 1 - Overwrite

Ghi đè dữ liệu cũ bằng dữ liệu mới:

```
1 -- Example: Update store address
2 UPDATE store_dim
3 SET address = 'New Address',
4     city = 'New City'
5 WHERE store_id = 123;
```

4.6.2 SCD Type 2 - Add New Row

Thêm bản ghi mới và giữ lại lịch sử:

```
1 def process_scd_type2(df, dim_table, key_col, attributes, conn):
2     """
3     Handle Slowly Changing Dimension Type 2
4     """
5     # Create a copy to avoid SettingWithCopyWarning
6     df = df.copy()
7
8     if current_dim.empty:
9         # Initial load - all records are new
10        df = df.assign(
11            start_date=datetime.now().date(),
12            end_date=None,
13            is_active=True
14        )
15        df.to_sql(dim_table, conn, if_exists='append', index=False)
16        return
17
18    # Find new and changed records
19    merged = df.merge(current_dim, on=key_col, how='outer',
20                      suffixes=('', '_current'), indicator=True)
21
22    # New records
23    new_records = merged[merged['_merge'] == 'left_only'][key_col + attributes].copy()
24    if not new_records.empty:
25        new_records = new_records.assign(
26            start_date=datetime.now().date(),
27            end_date=None,
28            is_active=True
29        )
30        new_records.to_sql(dim_table, conn, if_exists='append', index=False)
31
32    # Changed records
33    changed_records = merged[merged['_merge'] == 'both'].copy()
34    for attr in attributes:
35        changed_records = changed_records[
36            changed_records[attr] != changed_records[f'{attr}_current']
37        ]
38
39    if not changed_records.empty:
40        # Expire old record versions
41        conn.execute(text(f"""
42            UPDATE {dim_table}
43            SET end_date = :end_date, is_active = 0
44            WHERE {key_col} = :key_val AND is_active = 1
```



```

45         """), {'end_date': datetime.now().date(),
46               'key_val': row[key_col]})
47
48     # Insert new versions
49     new_versions = changed_records[key_col + attributes].copy()
50     new_versions = new_versions.assign(
51         start_date=datetime.now().date(),
52         end_date=None,
53         is_active=True
54     )
55     new_versions.to_sql(dim_table, conn, if_exists='append', index=False)

```

5. Advanced ETL Concepts

5.1 Change Data Capture (CDC)

CDC theo dõi và ghi lại mọi thay đổi trong database:

```

1 def implement_cdc_tracking():
2     """
3     Implement CDC tracking with timestamp and checksum
4     """
5     # Add CDC columns to the staging table
6     cdc_columns = {
7         'cdc_timestamp': 'TIMESTAMP DEFAULT CURRENT_TIMESTAMP',
8         'cdc_operation': 'VARCHAR(10)', # INSERT, UPDATE, DELETE
9         'cdc_checksum': 'VARCHAR(64)' # MD5 hash of the record
10    }
11
12    # Trigger to automatically record CDC
13    trigger_sql = """
14    CREATE TRIGGER cdc_trigger_sales
15    AFTER INSERT OR UPDATE OR DELETE ON sales_source
16    FOR EACH ROW
17    BEGIN
18        IF TG_OP = 'INSERT' THEN
19            INSERT INTO cdc_log VALUES (NEW.*, 'INSERT', CURRENT_TIMESTAMP);
20        ELSIF TG_OP = 'UPDATE' THEN
21            INSERT INTO cdc_log VALUES (NEW.*, 'UPDATE', CURRENT_TIMESTAMP);
22        ELSIF TG_OP = 'DELETE' THEN
23            INSERT INTO cdc_log VALUES (OLD.*, 'DELETE', CURRENT_TIMESTAMP);
24        END IF;
25    END;
26    """

```

5.2 Data Quality Monitoring

```

1 def data_quality_checks(df):
2     """
3     Check data quality
4     """
5     quality_report = {}
6
7     # 1. Completeness Check
8     quality_report['completeness'] = {
9         'total_records': len(df),
10        'missing_values': df.isnull().sum().to_dict(),

```

```

11     'completeness_rate': ((len(df) - df.isnull().sum()) / len(df) * 100).to_dict()
12 }
13
14 # 2. Uniqueness Check
15 duplicate_rate = df.duplicated().sum() / len(df) * 100
16 quality_report['uniqueness'] = {
17     'duplicate_records': df.duplicated().sum(),
18     'duplicate_rate': duplicate_rate
19 }
20
21 # 3. Validity Check
22 numeric_cols = df.select_dtypes(include=['number']).columns
23 invalid_numeric = (df[numeric_cols] < 0).sum()
24 quality_report['validity'] = {
25     'invalid_numeric_values': invalid_numeric.to_dict()
26 }
27
28 # 4. Consistency Check
29 business_rules_violations = []
30 if 'profit' in df.columns and 'revenue' in df.columns:
31     invalid_profit = (df['profit'] > df['revenue']).sum()
32     business_rules_violations.append({
33         'rule': 'profit_should_not_exceed_revenue',
34         'violations': invalid_profit
35     })
36
37 quality_report['consistency'] = business_rules_violations
38
39 return quality_report

```

5.3 Performance Optimization

5.3.1 Parallel Processing

```

1 from concurrent.futures import ThreadPoolExecutor, ProcessPoolExecutor
2 import multiprocessing as mp
3
4 def parallel_etl_processing(file_list, num_workers=None):
5     """
6     Parallel ETL processing for multiple files
7     """
8     if num_workers is None:
9         num_workers = mp.cpu_count()
10
11     def process_single_file(file_path):
12         try:
13             # Extract
14             df = pd.read_csv(file_path)
15
16             # Transform
17             df_transformed = transform_data(df)
18
19             # Load
20             load_to_warehouse(df_transformed)
21
22             return {"file": file_path, "status": "success", "records": len(df)}
23         except Exception as e:
24             return {"file": file_path, "status": "error", "error": str(e)}
25

```

```

26     # Use ProcessPoolExecutor for CPU-intensive tasks
27     with ProcessPoolExecutor(max_workers=num_workers) as executor:
28         results = list(executor.map(process_single_file, file_list))
29
30     return results

```

5.3.2 Memory Optimization

```

1 def memory_efficient_processing(file_path, chunk_size=50000):
2     """
3     Process large files with memory optimization
4     """
5     # Use categorical data type for string columns
6     dtype_map = {
7         'store': 'category',
8         'name': 'category',
9         'city': 'category',
10        'county': 'category',
11        'category_name': 'category',
12        'vendor_name': 'category'
13    }
14
15    # Read only necessary columns
16    required_columns = [
17        'invoice_line_no', 'date', 'store', 'sale_dollars',
18        'sale_bottles', 'state_bottle_cost', 'state_bottle_retail'
19    ]
20
21    chunk_iter = pd.read_csv(
22        file_path,
23        chunksize=chunk_size,
24        dtype=dtype_map,
25        usecols=required_columns,
26        parse_dates=['date']
27    )
28
29    for chunk in chunk_iter:
30        # Process each chunk
31        processed_chunk = transform_chunk(chunk)
32
33        # Load immediately to free memory
34        load_chunk_to_warehouse(processed_chunk)
35
36        # Explicitly delete to free memory
37        del chunk, processed_chunk

```

5.4 Error Handling và Monitoring

```

1 import logging
2 from datetime import datetime
3
4 def setup_etl_logging():
5     """
6     Setup logging for ETL pipeline
7     """
8     logging.basicConfig(
9         level=logging.INFO,

```

```

10         format='%(asctime)s - %(name)s - %(levelname)s - %(message)s',
11         handlers=[
12             logging.FileHandler(f'etl_pipeline_{datetime.now().strftime("%Y%m%d")}.log
13         '),
14             logging.StreamHandler()
15         ]
16     )
17     return logging.getLogger('ETL_Pipeline')
18
19 def robust_etl_pipeline(source_config, transform_config, target_config):
20     """
21     ETL pipeline with robust error handling
22     """
23     logger = setup_etl_logging()
24     pipeline_stats = {
25         'start_time': datetime.now(),
26         'processed_records': 0,
27         'error_records': 0,
28         'status': 'RUNNING'
29     }
30
31     try:
32         logger.info("Starting ETL Pipeline")
33
34         # Extract phase
35         logger.info("Starting extraction phase")
36         raw_data = extract_data(source_config)
37         logger.info(f"Extracted {len(raw_data)} records")
38
39         # Transform phase
40         logger.info("Starting transformation phase")
41         transformed_data, error_data = transform_data_with_error_handling(
42             raw_data, transform_config
43         )
44
45         pipeline_stats['processed_records'] = len(transformed_data)
46         pipeline_stats['error_records'] = len(error_data)
47
48         if len(error_data) > 0:
49             logger.warning(f"Found {len(error_data)} error records")
50             # Save error records for review
51             save_error_records(error_data)
52
53         # Load phase
54         logger.info("Starting load phase")
55         load_data(transformed_data, target_config)
56
57         pipeline_stats['status'] = 'SUCCESS'
58         logger.info("ETL Pipeline completed successfully")
59
60     except Exception as e:
61         pipeline_stats['status'] = 'FAILED'
62         logger.error(f"ETL Pipeline failed: {str(e)}")
63         raise
64
65     finally:
66         pipeline_stats['end_time'] = datetime.now()
67         pipeline_stats['duration'] = (
68             pipeline_stats['end_time'] - pipeline_stats['start_time']
69         ).total_seconds()

```

```
69         logger.info(f"Pipeline statistics: {pipeline_stats}")
70
71
72     return pipeline_stats
```