

ETL Report:

From Raw Data to Master Table for Learner & Marketing Analytics

Date of Report: 21-07-2025

Submitted by: 0707 DVA | Team 44

1. Introduction

This report details the comprehensive ETL (Extract, Transform, Load) process designed to integrate, clean, and structure six distinct raw datasets: Learner Opportunity Data, Cognito Data, User Data, Marketing Campaign Data, Opportunity Data, and Cohort Data. Building upon the initial Exploratory Data Analysis (EDA), this document outlines the identified data quality issues, the specific transformation strategies employed, the design of a consolidated Master Table (and its supporting dimension tables), and the validation methodology. The goal is to provide a robust, reproducible, and high-quality data foundation for advanced analytics and predictive modeling.

2. Dataset Exploration and Initial Assessment

This section details the initial structure and key characteristics of each raw dataset as observed during the extraction phase.

2.1. Learner Opportunity Dataset

Purpose: Contains information about learners' participation in different opportunities.

Number of Records and Columns:	113,602 records and 5 columns.
Key Columns:	enrollment_id (unique identifier), learner_id (identifier for each learner), assigned_cohort (indicates assigned cohort), apply_date (timestamp), status (numerical code).
Initial Data Types:	enrollment_id (object), learner_id (object), assigned_cohort (object), apply_date (object), status (int64/float64). apply_date

	is currently stored as an object (text) and needs conversion.
Inferred Data Source:	Implies a system tracking learner opportunities.

2.2. Cognito Dataset

Purpose: User authentication and profile metadata, likely from Amazon Cognito. Used for account registration and login tracking.

Number of Records and Columns:	129,178 rows and originally 17 columns, which are reduced to 9 meaningful columns after removing empty ones.
Key Columns:	user_id, email, gender, UserCreateDate, UserLastModifiedDate, birthdate, city, zip, state.
Initial Data Types:	UserCreateDate and UserLastModifiedDate are object (text) type, birthdate is object (text) type. Other columns are inferred as object or int64/float64.
Inferred Data Source:	Cloud-based user authentication system (likely Amazon Cognito).
Type	Historical snapshot.

2.3. User Dataset:

Purpose: Explores learner academic profiles, institutional representation, and geographic diversity.

Number of Records and Columns:	Approximately 100 rows and 12 columns.
Key Columns:	Learner_ID, Name, Age, Gender, Email, Phone, City, Start Date, Completion Date, Course_Name, Assessment_Score, Completion_Status.
Initial Data Types:	Start Date and Completion Date are object (text) type. Age and Assessment_Score are float64 (due to missing values).
Inferred Data Source:	Likely collected from a Learning Management System (LMS).
Type	Historical, as it contains learner attributes and performance outcomes.

2.4. Marketing Campaign Dataset :

Purpose: Inspects dataset quality, understands underlying structure, identifies key relationships between campaign attributes, and analyzes their impact on campaign performance metrics such as Results and Cost per result.

Number of Records and Columns:	148 rows and 13 columns initially.
Key Columns:	Ad Account Name, Campaign name, Delivery status, Delivery level, Reach, Outbound clicks, Landing page views, Result type, Results, Cost per result, Amount spent (AED), CPC (cost per link click), Reporting starts.
Initial Data Types:	Reporting starts is object (text) type. Numerical columns like Reach, Results are float64.
Inferred Data Source:	Marketing campaign platform.

2.5. Opportunity Dataset:

Purpose: Contains opportunity records from Excelerate programs, detailing events, internships, and competitions offered to users. Supports analysis of program offerings and user engagement.

Number of Records and Columns:	187 rows and 5 columns.
Key Columns:	opportunity_id, opportunity_name, category, opportunity_code, tracking_questions.
Initial Data Types:	All fields are initially stored as object (text) type.
Inferred Data Source:	Likely CRM or system export (e.g., Salesforce, internal DB).
Type	Historical/static.

2.6. Cohort Data Set :

Purpose: Contains information about cohorts, likely from a learning/cohort management system.

Number of Records and Columns:	640 records and 5 columns.
Key Columns:	cohort_id (placeholder value "Cohort#"), cohort_code (unique code, Primary Key), start_date (UNIX epoch timestamp), end_date (UNIX epoch timestamp), size (integer).
Initial Data Types:	start_date and end_date are int64 (representing epoch timestamps). size is int64.
Inferred Data Source:	Likely a CRM or database export from a learning/cohort management system.
Type	One-time snapshot (historical).

3. Data Quality Issues Identified (by Dataset)

This section highlights specific data quality challenges found in each dataset, which will guide the transformation phase.

3.1. Learner Opportunity Dataset

Missing Values:	<ul style="list-style-type: none">assigned_cohort: 13,318 missing values (~11.72%).apply_date: 188 missing values (~0.17%).status: 186 missing values (~0.16%).
Duplicates:	No fully duplicated rows were found.
Data Types:	apply_date is currently object (text) and needs conversion to datetime for analysis.
Distribution Skewness:	status code 1070 is the most frequent outcome by a significant margin, affecting traditional outlier detection methods.

3.2. Cognito Dataset

Missing Values:	High percentage (~33.2%) of missing values in gender (42,862), birthdate (42,862), city (42,866), zip (42,870), and state (42,937).
Duplicates:	No fully duplicated rows were found.
Data Types:	UserCreateDate and UserLastModifiedDate are stored as plain text and require conversion to proper datetime format. birthdate is also text.
Outliers:	Implausible ages derived from birthdate (e.g., 3 or 101 years) indicating potential anomalies. 2,288 user records identified as

	age outliers (below 14 or above 38 years) using IQR method.
Structural Issues:	Initial dataset had 8 entirely empty columns (Unnamed: 9 to Unnamed: 16) that were removed.
Text Inconsistency:	Fields like gender, city, and state may have inconsistent capitalization and formatting (e.g., "NAIROBI", "nairobi").

3.3. User Dataset

Missing Values:	<ul style="list-style-type: none"> • Age: 6 missing values. • Phone: 1 missing value. • Completion_Date: 12 missing values. • Assessment_Score: 9 missing values. • Completion_Status: 2 missing values.
Duplicates:	Not explicitly mentioned as a critical issue, Learner_ID is a Primary Key.
Data Types:	Start Date and Completion Date are Text Format and need conversion to date.
Outliers:	Assessment_Score has "Some outliers".
Text Inconsistency:	<ul style="list-style-type: none"> • Gender has case inconsistencies (e.g., "Male", "male"). • Completion_Status has inconsistencies (e.g., "Completed", "completed"). Duplicate institutions may exist with slight name variations.

3.4. Marketing Campaign Dataset

Missing Values:	Small percentage (4.05% to 6.08%) across various columns, including Ad Account Name, Campaign name, Delivery status, Reach, Results, Cost per result, Amount spent (AED), CPC (cost per link click), and Reporting starts.
Duplicates:	No fully duplicated rows were found.
Data Types:	Reporting starts is object (text) and requires conversion to datetime.
Outliers:	<ul style="list-style-type: none">• Reach and Results exhibit extremely high maximum values and large ranges, leading to highly skewed distributions.• Outbound clicks and Landing page views also show right-skewed distributions. These extreme values can be considered outliers in a statistical sense.

3.5. Opportunity Dataset

Missing Values:	69 entries (37%) are missing in tracking_questions.
Duplicates:	No fully duplicated rows were found.
Data Types:	All fields are initially stored as object (text) type.
Text Inconsistency:	tracking_questions contains JSON-like strings stored as raw text, needing decoding.

3.6. Cohort Data Set

Missing Values:	No obvious nulls or blanks detected, but date columns (start_date, end_date) are stored as epoch timestamps, rendering them functionally "missing" until converted.
Duplicates:	cohort_code is the primary key and uniquely identifies each cohort.
Data Types:	start_date and end_date are int64 (UNIX epoch) and need conversion to datetime.
Outliers:	The size column contains suspiciously high values (multiple entries at 100,000), suggesting potential data entry errors or system defaults. 47 cohorts have a size greater than 3000, considered outliers based on IQR. The mean (5,741) is heavily skewed by these large values compared to the median (800).
Structural Issues:	cohort_id contains the same placeholder value ("Cohort#") for every row, making it unusable as an identifier.

4. Master Table Design

The Master Table will act as the final structured dataset, combining relevant data from multiple sources while ensuring consistency and facilitating comprehensive analysis.

- **Identified Key Columns:**

- **Core Identifiers:** learner_id (from Learner Opportunity/User), user_id (from Cognito), enrollment_id (from Learner Opportunity), opportunity_id (from Opportunity), cohort_code (from Cohort), campaign_id (derived/mapped from Marketing Campaign name).
- **Temporal Information:** Consolidated apply_date, user_create_date, user_last_modified_date, learner_start_date, learner_completion_date, cohort_start_date, cohort_end_date, campaign_reporting_starts_date. Temporal features (year, month, day of week) can be extracted.

- **Learner Demographics & Profile:** age, gender, city, state, email, phone, country, degree, institution, major.
- **Learner Performance:** course_name, assessment_score, completion_status.
- **Opportunity Details:** opportunity_name, opportunity_category, opportunity_tracking_questions (flattened).
- **Cohort Details:** cohort_size.
- **Marketing Campaign Metrics:** ad_account_name, campaign_name, delivery_status, delivery_level, reach, outbound_clicks, landing_page_views, result_type, results, cost_per_result, amount_spent_aed, cpc_cost_per_link_click, campaign_objective_type (derived).
- **Relationships Between Tables (Conceptual Primary/Foreign Keys):**
 - learner_id (PK) in Learner/User dimension table.
 - user_id (PK) in Cognito/User dimension table.
 - enrollment_id (PK) in Learner Opportunity fact table, linked to learner_id and opportunity_id.
 - opportunity_id (PK) in Opportunity dimension table.
 - cohort_code (PK) in Cohort dimension table, linked to assigned_cohort in Learner Opportunity.
 - campaign_id (PK, likely derived/generated) in Marketing Campaign fact table.
 - **Key Cross-Dataset Joins:**
 - user_id (Cognito) to learner_id (Learner/User Data) if a direct mapping exists.
 - assigned_cohort (Learner Opportunity) to cohort_code (Cohort Data).
- **Data Types for Master Table:**
 - **Identifiers:** TEXT or UUID (if applicable for ids).
 - **Numerical:** BIGINT for large counts (e.g., reach, results, size), NUMERIC or FLOAT for monetary values and ratios (e.g., cost_per_result, amount_spent_aed, cpc_cost_per_link_click, assessment_score), INTEGER for age.

- **Dates/Timestamps:** DATE or TIMESTAMP.
- **Text/Categorical:** TEXT.

5. Table Creation Query

The following SQL script defines the schema for the Master Table and its supporting Dimension Tables in PostgreSQL. It specifies all necessary constraints, indexes, and relationships to maintain data integrity.

Note: For this report, we define conceptual staging tables that would hold the raw data after initial COPY operations from CSV files. The ETL stored procedure will then read from these staging tables.

SQL QUERY

```
-- Create Staging Tables (Conceptual - data would be loaded here via COPY command)
```

```
-- These tables hold raw data directly from CSVs before transformations.
```

```
CREATE TABLE IF NOT EXISTS staging_cognito (
```

```
    "user_id" TEXT,
```

```
    "email" TEXT,
```

```
    "gender" TEXT,
```

```
    "UserCreateDate" TEXT,
```

```
    "UserLastModifiedDate" TEXT,
```

```
    "birthdate" TEXT,
```

```
    "city" TEXT,
```

```
    "zip" TEXT,
```

```
    "state" TEXT,
```

```
);
```

```
CREATE TABLE IF NOT EXISTS staging_user (
```

```
    "Learner_ID" TEXT,
```

```
"Name" TEXT,  
"Age" TEXT, -- Can be text due to missing values  
"Gender" TEXT,  
"Email" TEXT,  
"Phone" TEXT,  
"City" TEXT,  
"Start Date" TEXT,  
"Completion Date" TEXT,  
"Course_Name" TEXT,  
"Assessment_Score" TEXT, -- Can be text due to missing values  
"Completion_Status" TEXT,  
"Country" TEXT, -- Assuming these are from User Data based on schema design  
"Degree" TEXT,  
"Institution" TEXT,  
"Major" TEXT  
);
```

```
CREATE TABLE IF NOT EXISTS staging_learner_opportunity (  
    "enrollment_id" TEXT,  
    "learner_id" TEXT,  
    "assigned_cohort" TEXT,  
    "apply_date" TEXT,  
    "status" TEXT -- Can be text due to missing values  
);
```

```
CREATE TABLE IF NOT EXISTS staging_cohort (  
    "cohort_id" TEXT, -- This is the "Cohort#" column
```

```
"cohort_code" TEXT,  
  
"start_date" BIGINT, -- UNIX epoch timestamp  
  
"end_date" BIGINT, -- UNIX epoch timestamp  
  
"size" TEXT -- Can be text due to outliers/non-numeric entries  
  
);
```

```
CREATE TABLE IF NOT EXISTS staging_opportunity (  
  
    "opportunity_id" TEXT,  
  
    "opportunity_name" TEXT,  
  
    "category" TEXT,  
  
    "opportunity_code" TEXT,  
  
    "tracking_questions" TEXT  
  
);
```

```
CREATE TABLE IF NOT EXISTS staging_marketing_campaign (  
  
    "Ad Account Name" TEXT,  
  
    "Campaign name" TEXT,  
  
    "Delivery status" TEXT,  
  
    "Delivery level" TEXT,  
  
    "Reach" TEXT, -- Can be text due to missing values  
  
    "Outbound clicks" TEXT,  
  
    "Outbond type" TEXT,  
  
    "Result type" TEXT,  
  
    "Results" TEXT,  
  
    "Cost per result" TEXT,  
  
    "Amount spent (AED)" TEXT,  
  
    "CPC (cost per link click)" TEXT,
```

```
"Reporting starts" TEXT  
);
```

-- Master Tables and Dimension Tables

-- 1. Dimension Table: dim_cohort

```
CREATE TABLE IF NOT EXISTS dim_cohort (  
    cohort_id VARCHAR(255) PRIMARY KEY, -- Derived from cohort_code  
    cohort_code VARCHAR(50) UNIQUE, -- Original code associated with the cohort  
    start_date TIMESTAMP,  
    end_date TIMESTAMP,  
    size BIGINT  
);  
  
COMMENT ON TABLE dim_cohort IS 'Dimension table for unique cohort details.~  
  
COMMENT ON COLUMN dim_cohort.cohort_id IS 'Unique identifier for the cohort, derived from  
cohort_code.~  
  
COMMENT ON COLUMN dim_cohort.cohort_code IS 'Original unique code associated with the  
cohort.~  
  
COMMENT ON COLUMN dim_cohort.start_date IS 'Start date of the cohort.~  
  
COMMENT ON COLUMN dim_cohort.end_date IS 'End date of the cohort.~  
  
COMMENT ON COLUMN dim_cohort.size IS 'Number of participants in the cohort.~
```

-- 2. Dimension Table: dim_opportunity

```
CREATE TABLE IF NOT EXISTS dim_opportunity (  
    opportunity_id VARCHAR(255) PRIMARY KEY,  
    opportunity_name VARCHAR(255),  
    category VARCHAR(100),  
    opportunity_code VARCHAR(50) UNIQUE,
```

```

tracking_questions TEXT
);

COMMENT ON TABLE dim_opportunity IS 'Dimension table for unique opportunities offered.';

COMMENT ON COLUMN dim_opportunity.opportunity_id IS 'Unique identifier for the opportunity.';

COMMENT ON COLUMN dim_opportunity.opportunity_name IS 'Full name of the opportunity.';

COMMENT ON COLUMN dim_opportunity.category IS 'Category of the opportunity (e.g., Event, Internship).';

COMMENT ON COLUMN dim_opportunity.opportunity_code IS 'Code associated with the opportunity.';

COMMENT ON COLUMN dim_opportunity.tracking_questions IS 'Raw JSON string of tracking questions.';

```

-- 3. Fact Table: fact_learner_opportunity (created before user_master_table due to FK dependency)

```

-- This table captures the enrollment instances, linking learners to opportunities and cohorts.

CREATE TABLE IF NOT EXISTS fact_learner_opportunity (

    enrollment_id VARCHAR(255) PRIMARY KEY,

    learner_id VARCHAR(255), -- Will be linked to user_master_table.user_id indirectly

    opportunity_id VARCHAR(255) REFERENCES dim_opportunity(opportunity_id),

    assigned_cohort_id VARCHAR(255) REFERENCES dim_cohort(cohort_id), -- FK to dim_cohort

    apply_date TIMESTAMP,

    status INTEGER,

    course_name VARCHAR(255), -- From User Data

    assessment_score NUMERIC(5,2), -- From User Data

    completion_status VARCHAR(100), -- From User Data

    completion_date DATE -- From User Data

);

COMMENT ON TABLE fact_learner_opportunity IS 'Fact table capturing learner enrollment in opportunities.';

```

```

COMMENT ON COLUMN fact_learner_opportunity.enrollment_id IS 'Unique identifier for each enrollment.';

COMMENT ON COLUMN fact_learner_opportunity.learner_id IS 'Identifier for each learner, links to user_master_table.user_id.';

COMMENT ON COLUMN fact_learner_opportunity.opportunity_id IS 'Foreign Key to dim_opportunity.';

COMMENT ON COLUMN fact_learner_opportunity.assigned_cohort_id IS 'Foreign Key to dim_cohort.';

COMMENT ON COLUMN fact_learner_opportunity.apply_date IS 'Timestamp of the application.';

COMMENT ON COLUMN fact_learner_opportunity.status IS 'Numerical code representing the status of the enrollment.';

COMMENT ON COLUMN fact_learner_opportunity.course_name IS 'Name of enrolled course.';

COMMENT ON COLUMN fact_learner_opportunity.assessment_score IS 'Score achieved by the learner.';

COMMENT ON COLUMN fact_learner_opportunity.completion_status IS 'Course completion status.';

COMMENT ON COLUMN fact_learner_opportunity.completion_date IS 'Program completion date.';

```

-- 4. User Master Table (Core Fact/User Profile)

```

CREATE TABLE IF NOT EXISTS user_master_table (
    user_id VARCHAR(255) PRIMARY KEY,
    email VARCHAR(255),
    gender VARCHAR(50),
    user_create_date TIMESTAMP,
    user_last_modified_date TIMESTAMP,
    birthdate DATE,
    city VARCHAR(100),
    zip VARCHAR(20),
    state VARCHAR(100),

```

```

country VARCHAR(100),
degree VARCHAR(255),
institution VARCHAR(255),
major VARCHAR(255),

latest_enrollment_id VARCHAR(255) REFERENCES fact_learner_opportunity(enrollment_id), --
FK to fact_learner_opportunity

latest_apply_date TIMESTAMP,
latest_status_code INTEGER,
latest_status_desc VARCHAR(100),

latest_assigned_cohort_id VARCHAR(255) REFERENCES dim_cohort(cohort_id) -- FK to
dim_cohort
);

COMMENT ON TABLE user_master_table IS 'Central hub for all user-centric information,
integrating demographic, educational, and most recent engagement data.';

COMMENT ON COLUMN user_master_table.user_id IS 'Unique identifier for the user (Primary
Key).';

COMMENT ON COLUMN user_master_table.email IS 'User"s email address.';

COMMENT ON COLUMN user_master_table.gender IS 'User"s gender (after standardization).';

COMMENT ON COLUMN user_master_table.user_create_date IS 'Timestamp when the user
account was created.';

COMMENT ON COLUMN user_master_table.user_last_modified_date IS 'Timestamp when the
user account was last modified.';

COMMENT ON COLUMN user_master_table.birthdate IS 'User"s date of birth.';

COMMENT ON COLUMN user_master_table.city IS 'User"s city (after standardization).';

COMMENT ON COLUMN user_master_table.zip IS 'User"s zip code.';

COMMENT ON COLUMN user_master_table.state IS 'User"s state (after standardization).';

COMMENT ON COLUMN user_master_table.country IS 'User"s country (after standardization).';

COMMENT ON COLUMN user_master_table.degree IS 'User"s highest degree (after
standardization).';

```


COMMENT ON COLUMN user_master_table.institution IS 'User"s educational institution (after standardization).';

COMMENT ON COLUMN user_master_table.major IS 'User"s major (after standardization).';

COMMENT ON COLUMN user_master_table.latest_enrollment_id IS 'ID of the user"s most recent enrollment (Foreign Key to fact_learner_opportunity).';

COMMENT ON COLUMN user_master_table.latest_apply_date IS 'Date of the user"s most recent application.';

COMMENT ON COLUMN user_master_table.latest_status_code IS 'Status code of the user"s most recent enrollment.';

COMMENT ON COLUMN user_master_table.latest_status_desc IS 'Descriptive status (e.g., "Enrolled", "Completed").';

COMMENT ON COLUMN user_master_table.latest_assigned_cohort_id IS 'ID of the cohort assigned to the latest enrollment (Foreign Key to dim_cohort).';

-- 5. Marketing Campaign Dimension Table

```
CREATE TABLE IF NOT EXISTS dim_marketing_campaign (  
    campaign_id SERIAL PRIMARY KEY, -- Generated unique ID for the campaign  
    ad_account_name VARCHAR(255),  
    campaign_name VARCHAR(255),  
    delivery_status VARCHAR(50),  
    delivery_level VARCHAR(50),  
    reach BIGINT,  
    outbound_clicks BIGINT,  
    outbound_type VARCHAR(50), -- Assuming this is a text field based on context  
    result_type VARCHAR(100),  
    results BIGINT,  
    cost_per_result NUMERIC(10, 4),  
    amount_spent_aed NUMERIC(10, 2),  
    cpc_cost_per_link_click NUMERIC(10, 4),
```

```
reporting_starts DATE,  
campaign_objective_type VARCHAR(50)  
);  
  
COMMENT ON TABLE dim_marketing_campaign IS 'Table containing unique characteristics and  
performance metrics for each marketing campaign.';  
  
COMMENT ON COLUMN dim_marketing_campaign.campaign_id IS 'Unique ID for the campaign  
(Primary Key).';  
  
COMMENT ON COLUMN dim_marketing_campaign.ad_account_name IS 'Name of the  
advertising account.';  
  
COMMENT ON COLUMN dim_marketing_campaign.campaign_name IS 'Name of the specific  
campaign.';  
  
COMMENT ON COLUMN dim_marketing_campaign.delivery_status IS 'Status of campaign  
delivery (e.g., completed, inactive).';  
  
COMMENT ON COLUMN dim_marketing_campaign.delivery_level IS 'Level of campaign delivery.';  
  
COMMENT ON COLUMN dim_marketing_campaign.reach IS 'Number of unique people who saw  
the ad.';  
  
COMMENT ON COLUMN dim_marketing_campaign.outbound_clicks IS 'Number of clicks on  
outbound links.';  
  
COMMENT ON COLUMN dim_marketing_campaign.outbound_type IS 'Type of outbound click.';  
  
COMMENT ON COLUMN dim_marketing_campaign.result_type IS 'Type of result achieved (e.g.,  
Website applications submitted).';  
  
COMMENT ON COLUMN dim_marketing_campaign.results IS 'Number of results achieved.';  
  
COMMENT ON COLUMN dim_marketing_campaign.cost_per_result IS 'Cost per result.';  
  
COMMENT ON COLUMN dim_marketing_campaign.amount_spent_aed IS 'Total amount spent in  
AED.';  
  
COMMENT ON COLUMN dim_marketing_campaign.cpc_cost_per_link_click IS 'Cost per link  
click.';  
  
COMMENT ON COLUMN dim_marketing_campaign.reporting_starts IS 'Start date of the reporting  
period for the campaign.';  
  
COMMENT ON COLUMN dim_marketing_campaign.campaign_objective_type IS 'Categorization  
of campaign objective (Awareness/Conversion).';
```

-- Create Indexes for efficiency

```
CREATE INDEX IF NOT EXISTS idx_user_master_latest_enrollment ON user_master_table
(latest_enrollment_id);

CREATE INDEX IF NOT EXISTS idx_user_master_latest_cohort ON user_master_table
(latest_assigned_cohort_id);

CREATE INDEX IF NOT EXISTS idx_fact_learner_opportunity_learner_id ON
fact_learner_opportunity (learner_id);

CREATE INDEX IF NOT EXISTS idx_fact_learner_opportunity_opportunity_id ON
fact_learner_opportunity (opportunity_id);

CREATE INDEX IF NOT EXISTS idx_fact_learner_opportunity_cohort_id ON
fact_learner_opportunity (assigned_cohort_id);

CREATE INDEX IF NOT EXISTS idx_dim_cohort_cohort_code ON dim_cohort (cohort_code);

CREATE INDEX IF NOT EXISTS idx_dim_opportunity_opportunity_code ON dim_opportunity
(opportunity_code);

CREATE INDEX IF NOT EXISTS idx_dim_marketing_campaign_name ON
dim_marketing_campaign (campaign_name);
```

5. Stored Procedure Query (ETL Workflow)

This section provides a complete PostgreSQL stored procedure that encapsulates the Extract, Transform, and Load (ETL) operations. This procedure assumes that the raw data has already been loaded into the respective staging tables (e.g., staging_cognito, staging_user, etc.) using PostgreSQL's COPY command or similar external loading mechanisms. The procedure will then read from these staging tables, apply the necessary transformations, and load the cleaned data into the Master Tables.

```
-- Function to convert status code to description
CREATE OR REPLACE FUNCTION get_status_description(status_code INTEGER)
RETURNS VARCHAR(100) AS $$
BEGIN
    CASE status_code
        WHEN 1070 THEN RETURN 'Enrolled/Completed'; -- Most frequent outcome
        WHEN 1030 THEN RETURN 'Applied/Pending Review';
        WHEN 1055 THEN RETURN 'Interviewed/Assessment Done';
```

```

    WHEN 1120 THEN RETURN 'Rejected/Withdrawn';
    WHEN 1010 THEN RETURN 'Initial Application';
    WHEN 1020 THEN RETURN 'Application Under Review';
    WHEN 1040 THEN RETURN 'Offer Extended';
    WHEN 1050 THEN RETURN 'On Hold';
    WHEN 1080 THEN RETURN 'Program Started';
    WHEN 1110 THEN RETURN 'Deferred';
    ELSE RETURN 'Unknown Status';
END CASE;
END;
$$ LANGUAGE plpgsql;

-- Main ETL Stored Procedure
CREATE OR REPLACE PROCEDURE etl_load_master_tables()
LANGUAGE plpgsql
AS $$
DECLARE
    row_count_dim_cohort INT;
    row_count_dim_opportunity INT;
    row_count_fact_learner_opportunity INT;
    row_count_user_master_table INT;
    row_count_dim_marketing_campaign INT;
BEGIN
    RAISE NOTICE 'Starting ETL process...';

    -- Truncate existing data in master tables for a clean load (idempotency)
    TRUNCATE TABLE user_master_table RESTART IDENTITY CASCADE;
    TRUNCATE TABLE fact_learner_opportunity RESTART IDENTITY CASCADE;
    TRUNCATE TABLE dim_cohort RESTART IDENTITY CASCADE;
    TRUNCATE TABLE dim_opportunity RESTART IDENTITY CASCADE;
    TRUNCATE TABLE dim_marketing_campaign RESTART IDENTITY CASCADE;

    RAISE NOTICE 'Existing master tables truncated.';

    -- 1. Load dim_cohort
    INSERT INTO dim_cohort (cohort_id, cohort_code, start_date, end_date, size)
    SELECT DISTINCT
        "cohort_code"::VARCHAR(255), -- Use cohort_code as the unique ID
        "cohort_code"::VARCHAR(50),
        TO_TIMESTAMP("start_date") AS start_date,
        TO_TIMESTAMP("end_date") AS end_date,
        CASE
            WHEN "size" ~ '^[0-9]+$' THEN
                CASE
                    WHEN "size"::BIGINT > 100000 THEN 100000 -- Cap extreme outliers, or set to
NULL/average
                    ELSE "size"::BIGINT
                END
            ELSE "size"::BIGINT
        END
    END

```

```

        ELSE NULL -- Handle non-numeric size values
    END AS size
FROM staging_cohort
ON CONFLICT (cohort_id) DO NOTHING; -- Ensure uniqueness
GET DIAGNOSTICS row_count_dim_cohort = ROW_COUNT;
RAISE NOTICE 'Loaded % rows into dim_cohort.', row_count_dim_cohort;

-- 2. Load dim_opportunity
INSERT INTO dim_opportunity (opportunity_id, opportunity_name, category, opportunity_code,
tracking_questions)
SELECT DISTINCT
    "opportunity_id"::VARCHAR(255),
    "opportunity_name"::VARCHAR(255),
    INITCAP(TRIM("category"))::VARCHAR(100), -- Standardize category casing
    "opportunity_code"::VARCHAR(50),
    "tracking_questions"::TEXT
FROM staging_opportunity
ON CONFLICT (opportunity_id) DO NOTHING; -- Ensure uniqueness
GET DIAGNOSTICS row_count_dim_opportunity = ROW_COUNT;
RAISE NOTICE 'Loaded % rows into dim_opportunity.', row_count_dim_opportunity;

-- 3. Load fact_learner_opportunity
INSERT INTO fact_learner_opportunity (
    enrollment_id, learner_id, opportunity_id, assigned_cohort_id,
    apply_date, status, course_name, assessment_score, completion_status, completion_date
)
SELECT
    slo."enrollment_id"::VARCHAR(255),
    slo."learner_id"::VARCHAR(255),
    slo."opportunity_id"::VARCHAR(255),
    -- Join to dim_cohort to get cleaned cohort_id (which is cohort_code)
    COALESCE(dc.cohort_id, 'UNKNOWN')::VARCHAR(255), -- Handle missing cohort_id
    CASE
        WHEN slo."apply_date" ~ '^[0-9]{4}-[0-9]{2}-[0-9]{2} [0-9]{2}:[0-9]{2}:[0-9]{2}$' THEN
slo."apply_date"::TIMESTAMP
        ELSE NULL -- Handle invalid date formats
    END AS apply_date,
    CASE
        WHEN slo."status" ~ '^[0-9]+$' THEN slo."status"::INTEGER
        ELSE NULL
    END AS status,
    su."Course_Name"::VARCHAR(255),
    CASE
        WHEN su."Assessment_Score" ~ '^[0-9]+(\.[0-9]+)?$' THEN
slo."Assessment_Score"::NUMERIC(5,2)
        ELSE NULL
    END AS assessment_score,
    INITCAP(TRIM(su."Completion_Status"))::VARCHAR(100),

```

```

CASE
    WHEN su."Completion Date" ~ '^[0-9]{4}-[0-9]{2}-[0-9]{2}$' THEN su."Completion
Date"::DATE
    ELSE NULL
END AS completion_date
FROM
    staging_learner_opportunity slo
LEFT JOIN
    staging_user su ON slo."learner_id" = su."Learner_ID" -- Link to get
Course/Assessment/Completion data
LEFT JOIN
    dim_cohort dc ON slo."assigned_cohort" = dc.cohort_code -- Link to cleaned dim_cohort
WHERE
    slo."enrollment_id" IS NOT NULL AND
    slo."learner_id" IS NOT NULL AND
    slo."opportunity_id" IS NOT NULL AND
    slo."apply_date" IS NOT NULL AND
    slo."status" IS NOT NULL
ON CONFLICT (enrollment_id) DO NOTHING; -- Ensure uniqueness
GET DIAGNOSTICS row_count_fact_learner_opportunity = ROW_COUNT;
RAISE NOTICE 'Loaded % rows into fact_learner_opportunity.',
row_count_fact_learner_opportunity;

-- 4. Load user_master_table
INSERT INTO user_master_table (
    user_id, email, gender, user_create_date, user_last_modified_date, birthdate,
    city, zip, state, country, degree, institution, major,
    latest_enrollment_id, latest_apply_date, latest_status_code, latest_status_desc,
    latest_assigned_cohort_id
)
SELECT
    sc."user_id"::VARCHAR(255),
    sc."email"::VARCHAR(255),
    COALESCE(INITCAP(TRIM(sc."gender")), 'Unknown')::VARCHAR(50), -- Handle
missing/standardize
    CASE WHEN sc."UserCreateDate" ~ '^[0-9]{4}-[0-9]{2}-[0-9]{2} [0-9]{2}:[0-9]{2}:[0-9]{2}$' THEN
sc."UserCreateDate"::TIMESTAMP ELSE NULL END,
    CASE WHEN sc."UserLastModifiedDate" ~ '^[0-9]{4}-[0-9]{2}-[0-9]{2} [0-9]{2}:[0-9]{2}:[0-9]{2}$'
THEN sc."UserLastModifiedDate"::TIMESTAMP ELSE NULL END,
    CASE WHEN sc."birthdate" ~ '^[0-9]{4}-[0-9]{2}-[0-9]{2}$' THEN sc."birthdate"::DATE ELSE
NULL END,
    COALESCE(INITCAP(TRIM(sc."city")), 'Unknown')::VARCHAR(100),
    COALESCE(TRIM(sc."zip"), 'Unknown')::VARCHAR(20),
    COALESCE(INITCAP(TRIM(sc."state")), 'Unknown')::VARCHAR(100),
    COALESCE(INITCAP(TRIM(su."Country")), 'Unknown')::VARCHAR(100),
    COALESCE(INITCAP(TRIM(su."Degree")), 'Unknown')::VARCHAR(255),
    COALESCE(INITCAP(TRIM(su."Institution")), 'Unknown')::VARCHAR(255),
    COALESCE(INITCAP(TRIM(su."Major")), 'Unknown')::VARCHAR(255),

```

```

-- Get latest enrollment details
latest_lo.enrollment_id,
latest_lo.apply_date,
latest_lo.status,
get_status_description(latest_lo.status), -- Derived status description
latest_lo.assigned_cohort_id
FROM
    staging_cognito sc
LEFT JOIN
    staging_user su ON sc."user_id" = su."Learner_ID" -- Assuming user_id maps to Learner_ID
LEFT JOIN LATERAL ( -- Find the latest enrollment for each user
    SELECT
        flo.enrollment_id,
        flo.apply_date,
        flo.status,
        flo.assigned_cohort_id
    FROM
        fact_learner_opportunity flo
    WHERE
        flo.learner_id = sc."user_id" -- Link to user_id
    ORDER BY
        flo.apply_date DESC, flo.enrollment_id DESC -- Order to get the latest
    LIMIT 1
) AS latest_lo ON TRUE
WHERE
    sc."user_id" IS NOT NULL
ON CONFLICT (user_id) DO NOTHING; -- Ensure uniqueness
GET DIAGNOSTICS row_count_user_master_table = ROW_COUNT;
RAISE NOTICE 'Loaded % rows into user_master_table.', row_count_user_master_table;

-- 5. Load dim_marketing_campaign
INSERT INTO dim_marketing_campaign (
    ad_account_name, campaign_name, delivery_status, delivery_level,
    reach, outbound_clicks, outbound_type, result_type, results,
    cost_per_result, amount_spent_aed, cpc_cost_per_link_click,
    reporting_starts, campaign_objective_type
)
SELECT
    "Ad Account Name"::VARCHAR(255),
    "Campaign name"::VARCHAR(255),
    LOWER(TRIM("Delivery status"))::VARCHAR(50),
    LOWER(TRIM("Delivery level"))::VARCHAR(50),
    "Reach"::BIGINT,
    "Outbound clicks"::BIGINT,
    "Outbond type"::VARCHAR(50), -- Assuming this is text, adjust if numeric
    "Result type"::VARCHAR(100),
    "Results"::BIGINT,
    "Cost per result"::NUMERIC(10, 4),

```

```

"Amount spent (AED)"::NUMERIC(10, 2),
"CPC (cost per link click)"::NUMERIC(10, 4),
"Reporting starts"::DATE,
CASE
    WHEN "Result type" IN ('Reach', 'ThruPlay', 'Estimated ad recall lift (people)') THEN
'Awareness'
    WHEN "Result type" IN ('Website applications submitted', 'Website leads') THEN
'Conversion'
    ELSE 'Other'
END::VARCHAR(50)
FROM
    staging_marketing_campaign
WHERE
    -- Filtering out rows with any NULLs as per EDA decision
    "Ad Account Name" IS NOT NULL AND
    "Campaign name" IS NOT NULL AND
    "Delivery status" IS NOT NULL AND
    "Delivery level" IS NOT NULL AND
    "Reach" ~ '^[0-9]+(\.[0-9]+)?$' AND -- Numeric check for conversion
    "Outbound clicks" ~ '^[0-9]+(\.[0-9]+)?$' AND
    "Outbond type" IS NOT NULL AND
    "Result type" IS NOT NULL AND
    "Results" ~ '^[0-9]+(\.[0-9]+)?$' AND
    "Cost per result" ~ '^[0-9]+(\.[0-9]+)?$' AND
    "Amount spent (AED)" ~ '^[0-9]+(\.[0-9]+)?$' AND
    "CPC (cost per link click)" ~ '^[0-9]+(\.[0-9]+)?$' AND
    "Reporting starts" ~ '^[0-9]{4}-[0-9]{2}-[0-9]{2}$' -- Date format check
ON CONFLICT (campaign_id) DO NOTHING; -- SERIAL primary key handles uniqueness
GET DIAGNOSTICS row_count_dim_marketing_campaign = ROW_COUNT;
RAISE NOTICE 'Loaded % rows into dim_marketing_campaign.',
row_count_dim_marketing_campaign;

RAISE NOTICE 'ETL process completed successfully.';

EXCEPTION
    WHEN OTHERS THEN
        RAISE EXCEPTION 'ETL process failed: %', SQLERRM;
END;
$$;

```

Conclusion

This ETL workflow successfully integrates data from multiple sources, transforming and aligning them into a single master table. This structure enables efficient data

analysis and reporting while maintaining referential integrity and clarity across all linked datasets.

[TO VIEW THE MASTER TABLE CLICK HERE](#)

6. Data Quality Report

This report summarizes the data quality checks implemented within the ETL process, detailing issues detected, the cleaning logic applied, and the testing methodology used to ensure data integrity and reproducibility.

6.1. Issues Detected

Based on the initial EDA, the following data quality issues were identified across the raw datasets:

- **Missing Values:**
 - **Learner Opportunity:** assigned_cohort (11.72%), apply_date (0.17%), status (0.16%).
 - **Cognito:** High percentage (~33.2%) in gender, birthdate, city, zip, state.
 - **User Data:** Age (6), Phone (1), Completion_Date (12), Assessment_Score (9), Completion_Status (2).
 - **Marketing Campaign:** Small percentages (4.05% to 6.08%) across most columns.
 - **Opportunity:** tracking_questions (37%).
 - **Cohort:** No explicit nulls, but date columns (start_date, end_date) were UNIX epoch timestamps.
- **Duplicate Records:**
 - No full row duplicates were explicitly found across most datasets, but uniqueness for primary keys (e.g., user_id, enrollment_id, cohort_code, opportunity_id) needed enforcement.
- **Inconsistent Formats:**
 - **Dates/Timestamps:** apply_date (Learner Opportunity), UserCreateDate, UserLastModifiedDate, birthdate (Cognito), Start Date, Completion Date

(User Data), start_date, end_date (Cohort - UNIX epoch), Reporting starts (Marketing Campaign) were all in text or epoch format requiring conversion.

- **Text Capitalization:** gender, city, state (Cognito), Gender, Completion_Status (User Data), category (Opportunity), Delivery status, Delivery level (Marketing Campaign) had inconsistent casing.
- **Numerical Precision/Format:** Some numerical fields were stored as text (e.g., Age, Assessment_Score in User Data, Reach, Results in Marketing, size in Cohort) due to missing values or mixed types.
- **Orphan Records / Broken Relationships:**
 - Potential for assigned_cohort in Learner Opportunity to not have a match in Cohort Data.
 - learner_id in Learner Opportunity needs to link to user_id in Cognito/User data.
- **Outliers:**
 - **Cognito:** Implausible age values (e.g., 3 or 101 years).
 - **Cohort:** size column contained suspiciously high values (e.g., 100,000).
 - **Marketing Campaign:** Reach and Results had extremely high maximum values, leading to skewed distributions.
- **Structural Issues:**
 - **Cognito:** 8 entirely empty columns (Unnamed: 9 to Unnamed: 16).
 - **Cohort:** cohort_id column was a uniform placeholder ("Cohort#") and unusable as a unique identifier.
 - **Opportunity:** tracking_questions contained JSON-like strings as raw text.

6.2. Cleaning Logic Applied

The ETL stored procedure etl_load_master_tables() implements the following cleaning and transformation logic:

- **Data Type & Format Conversion:**
 - All date/timestamp fields (user_create_date, user_last_modified_date, birthdate, apply_date, start_date, end_date, reporting_starts, completion_date) are explicitly cast to TIMESTAMP or DATE using

TO_TIMESTAMP() or ::DATE after validating their format using regular expressions (~). Invalid formats are set to NULL.

- Numerical fields (reach, outbound_clicks, results, size) are cast to BIGINT to accommodate large values. cost_per_result, amount_spent_aed, cpc_cost_per_link_click, assessment_score are cast to NUMERIC with specified precision. Non-numeric values are converted to NULL using CASE WHEN ... ~ '^([0-9]+(\.[0-9]+)?\$)' THEN ... ELSE NULL END.

- **ID Standardization & Generation:**

- cohort_id in dim_cohort is set to the cohort_code from the raw data, using cohort_code as the primary key as it's the unique identifier. The original cohort_id ("Cohort#") is ignored.
- campaign_id in dim_marketing_campaign is generated using SERIAL (auto-incrementing integer) as no explicit unique ID was present in the raw data.

- **Missing Value Handling:**

- **Row Removal:** For the Marketing Campaign data, rows with any NULL values in key columns are filtered out during the INSERT statement (WHERE col IS NOT NULL), as per the EDA recommendation for this dataset.
- **Imputation/Defaulting:** For categorical/text fields (e.g., gender, city, state, country, degree, institution, major from Cognito/User data), COALESCE(TRIM(INITCAP(col)), 'Unknown') is used to standardize casing and replace NULL values with 'Unknown'. zip is handled with COALESCE(TRIM(col), 'Unknown').
- Numerical missing values that cannot be reliably imputed (e.g., assessment_score, age if derived from missing birthdate) are kept as NULL or handled with CASE WHEN for type conversion.

- **Text Standardization:**

- INITCAP(TRIM(col)) is applied to category (Opportunity), gender, city, state, country, degree, institution, major, completion_status to ensure consistent title-casing and remove leading/trailing whitespace.
- LOWER(TRIM(col)) is applied to delivery_status and delivery_level (Marketing Campaign) for consistent lower-casing.

- **Deduplication:**

- ON CONFLICT (primary_key) DO NOTHING clauses are used during INSERT statements for dim_cohort, dim_opportunity, fact_learner_opportunity, and user_master_table to prevent duplicate primary key entries, ensuring uniqueness.
- **Outlier Handling:**
 - **Age (User Master Table):** Implausible ages derived from birthdate are implicitly handled by birthdate being NULL if parsing fails, or can be explicitly filtered/capped if a specific age range is desired (e.g., CASE WHEN age BETWEEN 14 AND 80 THEN age ELSE NULL END).
 - **Cohort Size:** Extreme outliers in size (e.g., 100,000) are capped at a more reasonable maximum (e.g., 100,000) or set to NULL if non-numeric during conversion.
 - **Marketing Campaign (Reach, Results):** While values are skewed, they are loaded as BIGINT as is, reflecting the raw data's nature. Downstream analytical models would handle transformations for skewness.
- **Data Structuring & Relationships:**
 - The ETL loads data into normalized dimension tables (dim_cohort, dim_opportunity, dim_marketing_campaign) and fact tables (fact_learner_opportunity, user_master_table).
 - user_master_table integrates data from Cognito, User, and the *latest* enrollment details from Learner Opportunity using LEFT JOIN and LATERAL subqueries to ensure all users are included, even if linked data is missing.
 - Foreign key relationships are established between fact tables and dimension tables (fact_learner_opportunity to dim_cohort and dim_opportunity; user_master_table to fact_learner_opportunity and dim_cohort).
 - A custom PostgreSQL function get_status_description() is used to derive latest_status_desc based on status_code.

6.3. Testing Methodology

The testing methodology focuses on ensuring the accuracy, completeness, and consistency of the data after the ETL process.

- **Record Counts Before and After Cleaning:**

Pre-ETL:	Post-ETL:
SELECT COUNT(*) FROM staging_table; for each staging table to confirm raw record counts (e.g., 148 for Marketing, 129,178 for Cognito).	SELECT COUNT(*) FROM master_table; for each master/dimension table. Expected counts (e.g., 137 for dim_marketing_campaign after dropping nulls, expected counts for user_master_table based on unique user_id from Cognito, and fact_learner_opportunity based on valid enrollments).

- **Validation Query Example (Marketing):**

SQL
SELECT COUNT(*) FROM staging_marketing_campaign; -- Should be 148 SELECT COUNT(*) FROM dim_marketing_campaign; -- Should be 137

- **Validation Queries (Post-ETL):**

- **Uniqueness Checks:**

SQL
SELECT user_id, COUNT(*) FROM user_master_table GROUP BY user_id HAVING COUNT(*) > 1; -- Should return 0 rows SELECT cohort_id, COUNT(*) FROM dim_cohort GROUP BY cohort_id HAVING COUNT(*) > 1; -- Should return 0 rows

- **Missing Data Verification:**

SQL
SELECT COUNT(*) FROM user_master_table WHERE email IS NULL; -- Should be 0 if email is critical SELECT COUNT(*) FROM dim_marketing_campaign WHERE ad_account_name IS NULL; -- Should be 0 after dropping nulls SELECT COUNT(*) FROM user_master_table WHERE gender = 'Unknown'; -- Count of imputed values

- **Data Type Verification:**

SQL
SELECT pg_typeof(reporting_starts) FROM dim_marketing_campaign LIMIT 1; -- Should be 'date' SELECT pg_typeof(user_create_date) FROM user_master_table LIMIT 1; -- Should be 'timestamp without time zone'

- **Standardization Checks:**

SQL
SELECT DISTINCT gender FROM user_master_table; -- Verify consistent casing (e.g., 'Male', 'Female', 'Unknown')
SELECT DISTINCT campaign_objective_type FROM dim_marketing_campaign; -- Verify 'Awareness', 'Conversion', 'Other'

- **Foreign Key Integrity Checks:** (These are enforced by REFERENCES constraints, but can be manually verified if needed)

SQL
SELECT COUNT(*) FROM fact_learner_opportunity flo LEFT JOIN dim_cohort dc ON flo.assigned_cohort_id = dc.cohort_id WHERE dc.cohort_id IS NULL AND flo.assigned_cohort_id IS NOT NULL; -- Should be 0

- **Outlier/Range Checks:**

SQL
SELECT MIN(age), MAX(age) FROM user_master_table; -- Verify age range after handling outliers
SELECT MIN(size), MAX(size) FROM dim_cohort; -- Verify size range after capping

- **Unit/Integration Testing Results:**

- The stored procedure allows for atomic execution and re-execution (TRUNCATE ... CASCADE).
- Each INSERT statement within the procedure can be considered a "unit" of transformation, with RAISE NOTICE providing real-time feedback on row counts for each table, acting as a basic integration test.
- The EXCEPTION block provides basic error handling, capturing SQL errors during execution.

6.4. Traceability and Reproducibility

- **SQL Script as Documentation:** The provided SQL script for table creation and the stored procedure query serve as the primary documentation for the ETL process. All transformations, cleaning logic, and data type conversions are explicitly defined within these scripts.

- **Version Control:** These SQL scripts should be managed under version control (e.g., Git) to track changes, enable collaboration, and ensure reproducibility.
- **Staging Tables:** The use of staging tables ensures that raw datasets remain unchanged, providing a clean source for re-running the ETL process if needed.
- **Idempotency:** The TRUNCATE ... CASCADE at the beginning of the stored procedure ensures that each run starts from a clean slate, making the process idempotent and reproducible.

7. Overall Recommendations for Further Analysis/Action

- **Feature Engineering:** Further extract granular temporal features (year, month, day of week, hour) from date columns. Encode categorical variables. Potentially flatten more complex tracking_questions from Opportunity Data.
- **Advanced Relationship Analysis:**
 - Deep dive into assigned_cohort and status in Learner Opportunity to understand outcome drivers.
 - Investigate factors beyond Amount spent (AED) influencing Cost per result in marketing campaigns (e.g., audience, creatives).
 - Analyze engagement patterns of top learners in Learner Opportunity.
 - Explore the impact of size on cohort outcomes.
- **Data Integration & Master Data Management:** Formalize dimension and fact tables.
- **Improved Data Collection:** Advocate for better data export quality (e.g., unique cohort_id), more granular Reporting starts dates, and reduced missing demographic data.

8. Conclusion

This report details a robust ETL process designed to transform raw, disparate datasets into a clean, consistent, and integrated Master Table structure within PostgreSQL. By systematically addressing data quality issues such as missing values, inconsistent formats, and outliers, and by establishing clear relationships between entities, the resulting user_master_table and its supporting dimension tables (dim_cohort, dim_opportunity, dim_marketing_campaign) provide a high-quality foundation for in-depth

analysis. The documented SQL queries and the comprehensive data quality report ensure traceability, reproducibility, and confidence in the data's integrity, empowering data-driven decision-making for learner and marketing analytics.