

WEEK - 4

COMPILATION OF FINDINGS & FINAL PRESENTATION

DATA VISUALIZATION ASSOCIATE EARLY INTERNSHIP



PRESENTED BY

SLU 0704 | TEAM 15

MAY 2025

AGENDA

- Executive Summary
- Project Objective
- Business Problem
- Data Sources Overview
- ETL Process Overview
- Data Quality Improvements
- Master Table Creation
- Wireframe Model Design
- Key Insights
- Additional Insights
- Advanced Data Validation Metrics
- Challenges vs. Resolutions
- Impact Summary
- Next Steps & Final Takeaways

EXECUTIVE SUMMARY

Over the course of this project, six raw datasets were meticulously cleaned, standardized, and transformed into analytics-ready formats, culminating in the integration of five datasets into a robust Master Table, while the Marketing Campaign Data was analyzed independently due to structural limitations. Through the automation of ETL processes using advanced stored procedures, the overall data quality was significantly improved, setting the foundation for meaningful dashboard development and strategic decision-making.

PROJECT OBJECTIVE

The primary objective of this project was to systematically transform fragmented and inconsistent raw data into a structured, high-quality, and integrated dataset, ensuring analytics readiness. This transformation aimed to resolve critical issues such as data inconsistency, duplication, missing information, and system isolation by establishing a scalable, automated, and repeatable ETL process that would deliver trustworthy insights to drive strategic business initiatives.

BUSINESS PROBLEM

Challenges Identified in Raw Data -

- A significant volume of missing values, with up to 33% in some datasets.
- Textual and date inconsistencies affecting data reliability.
- No direct integration or relational mapping across datasets.
- Presence of duplicate and unreliable records, compromising trustworthiness.

Business Need -

The organization required a reliable, clean, and fully integrated data ecosystem to support effective reporting, dashboard creation, and informed strategic decision-making.

DATA SOURCES OVERVIEW

Datasets Processed -

Learner Data -

- Records - 129,259, Fields - 5
- Primary Key - learner_id

Opportunity Data -

- Records - 187, Fields - 5
- Primary Key - opportunity_id

Cohort Data -

- Records - 639, Fields - 5
- Primary Key - cohort_code

Learner-Opportunity Data -

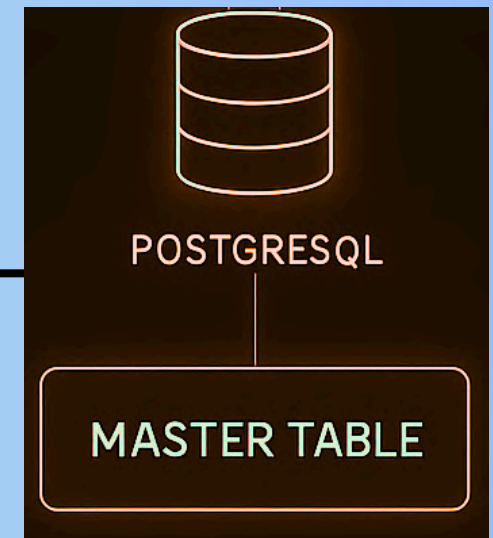
- Records - 113,602, Fields - 5
- Primary Key - enrollment_id

Cognito Data -

- Records - 129,178, Fields - 9
- Primary Key - user_id

Marketing Data (analyzed separately) -

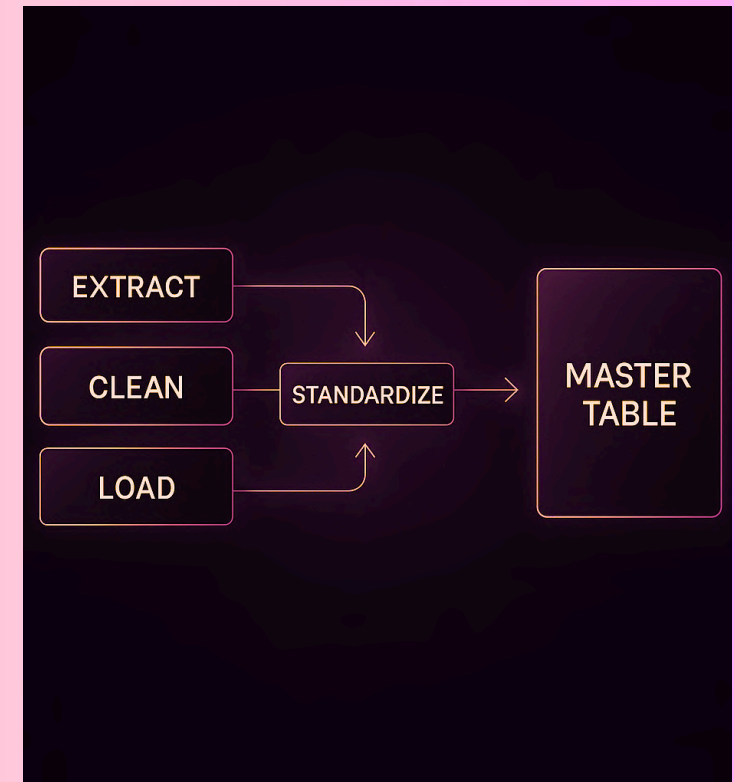
- Records - 155, Fields - 13
- No Unique Primary Key



ETL PROCESS OVERVIEW

Major Steps Implemented -

- Extraction of raw CSV files and systematic importation into PostgreSQL databases.
- Comprehensive cleaning of text fields, date formats, and numeric values using built-in SQL functions such as TRIM, INITCAP, UPPER, and TO_DATE.
- Standardization and normalization of key fields to ensure consistency across datasets.
- Loading and integration of cleaned data into master tables to form a consolidated data environment.



DATA QUALITY IMPROVEMENTS

- Removed unnecessary whitespace with TRIM().
- Standardized text casing using INITCAP() for better readability and analysis.
- Enforced upper casing for opportunity codes to maintain uniformity.
- Applied Regex Validation to ensure all date fields conformed to valid formats.
- Utilized DISTINCT selection to eliminate duplicate records without losing unique entries.
- Preserved NULL values in critical fields to maintain data authenticity.

MASTER TABLE CREATION

Integration Strategy -

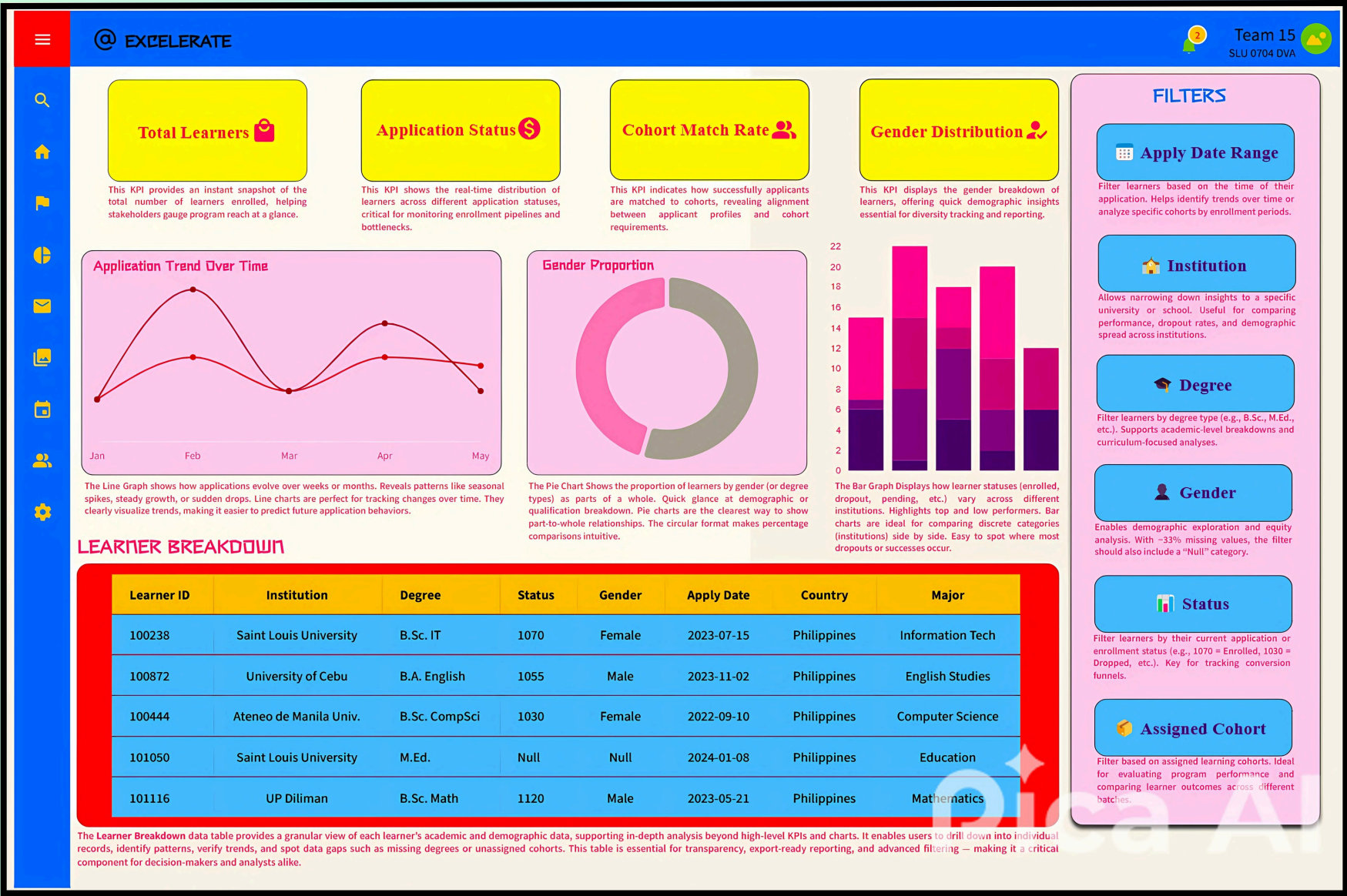
- Effective JOIN operations were performed across learner, opportunity, cohort, and user datasets, ensuring referential integrity.
- Marketing data was strategically excluded due to the absence of relatable foreign keys, leading to a separate analytical track.

Fields Used -

- enrollment_id, assigned_cohort, apply_date, status (learner_opportunity)
- learner_id, degree, institution, major, country (learner_raw)
- start_date, end_date, size (cohort)
- email, gender, birthday, city, zip, state (cognito)
- opportunity_id, opportunity_name, category, tracking_questions (opportunity_raw)

Master table integrates 5 out of 6 datasets. Marketing Dataset is excluded due to structural limitations. Data quality is acceptable for analysis; improvements noted. The pipeline is automated using a stored procedure and ready for dashboard use.

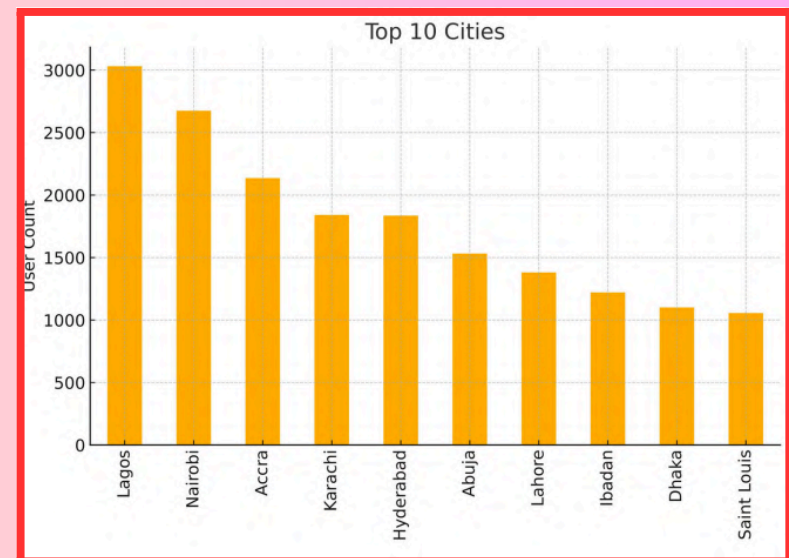
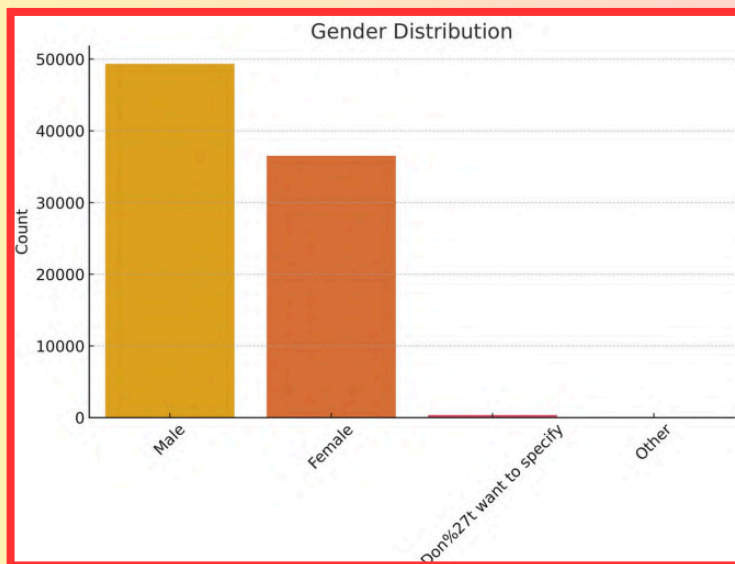
WIREFRAME MODEL DESIGN



KEY INSIGHT #1 - MISSING DEMOGRAPHIC DATA

Approximately 33% of entries in the Cognito dataset lacked critical demographic fields such as gender, city, and state, posing challenges to personalization and segmentation strategies.

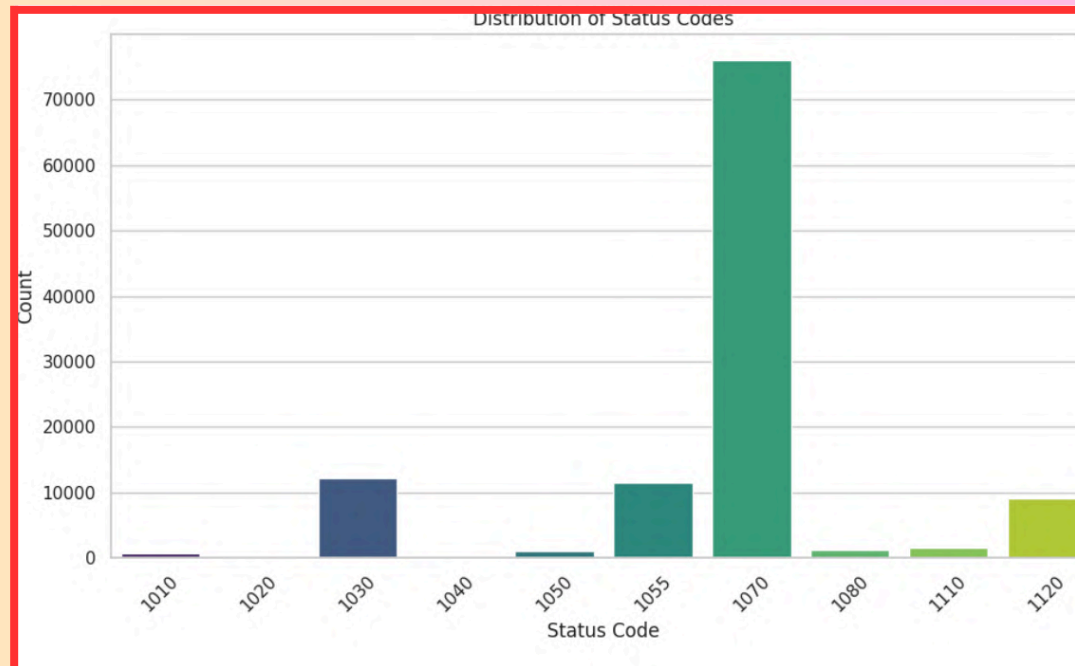
- **Impact** - The incomplete data limited the organization's ability to tailor experiences and communication.
- **Recommendation** - Strengthen data collection practices during user registration to improve demographic completeness.



KEY INSIGHT #2 - ENROLLMENT STATUS DOMINANCE

Enrollment Status Code 1070 accounted for nearly 67% of all records, suggesting either a program bottleneck or a reporting milestone issue.

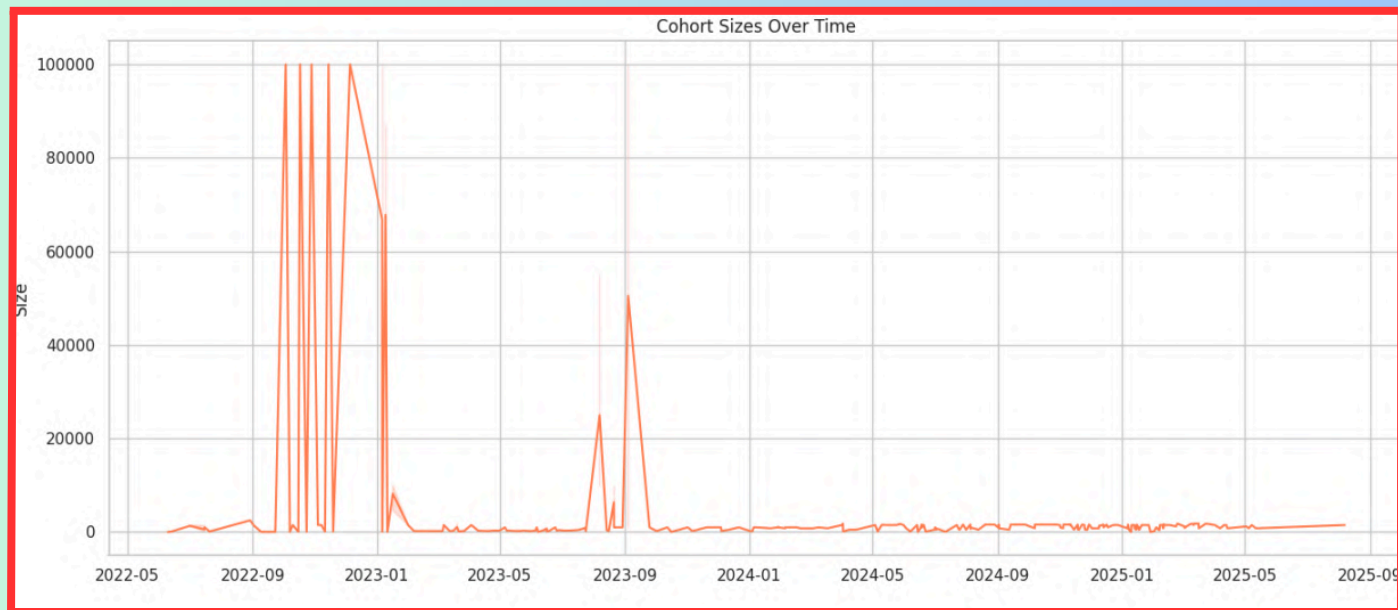
- **Impact** - This dominance risks masking underlying program engagement issues.
- **Recommendation** - Conduct a deep-dive analysis into the lifecycle of Status Code 1070 to uncover root causes.



KEY INSIGHT #3 - COHORT ANOMALIES

Significant variation was observed in cohort sizes and durations, with some cohorts showing 0-day durations, highlighting potential data entry or program structure issues.

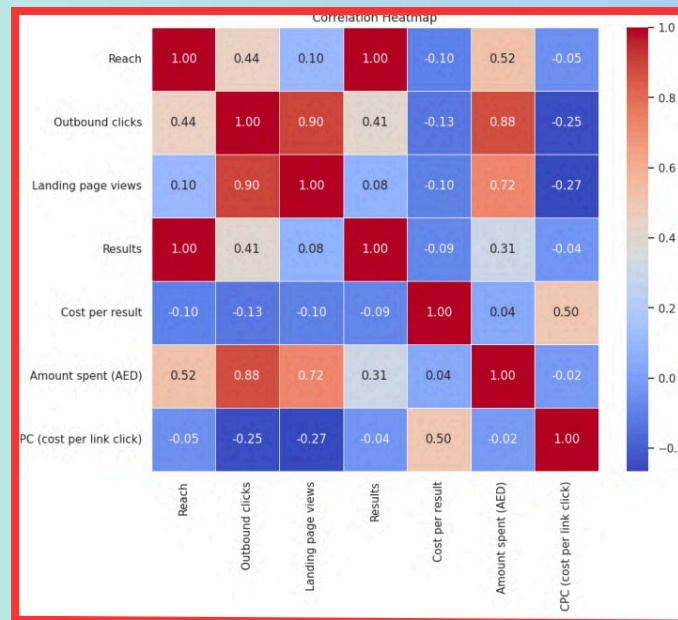
- **Impact** - Unreliable cohort data complicates capacity planning and program evaluation.
- **Recommendation** - Implement strict data validation rules at the point of cohort creation.



KEY INSIGHT #4 - MARKETING DATA ISOLATION

The Marketing Campaign Data lacked any joinable learner identifiers, preventing meaningful attribution of campaign performance to learner outcomes.

- **Impact** - Campaign ROI analyses were disconnected from learner enrollment or progression data.
- **Recommendation** - Introduce UTM tagging and tracking IDs in marketing initiatives to ensure traceability.



ADDITIONAL INSIGHTS

#1 OUTLIER DETECTION -

Outlier detection revealed users with improbable birthdates, some suggesting ages below 5 years or above 100 years, indicating flawed input validation.

- **Impact** - Such anomalies skew demographic and segmentation analyses, leading to unreliable insights.
- **Recommendation** - Introduce robust front-end validations to ensure realistic and accurate user data.

#2 LEARNER ENGAGEMENT CONCENTRATION -

Despite having over 129,000 learners in the system (from Learner_Raw), only 187 unique learners were responsible for generating 113,602 enrollment records in LearnerOpportunity_Raw. This means that 0.14% of all learners accounted for almost all opportunity engagement!

- **Impact** - The broader learner population is either inactive, undocumented, or engaging outside tracked opportunities.
- **Recommendation** - Expand efforts to activate the remaining 99.86% of users through targeted outreach, nudges, or product redesign.








ADVANCED DATA VALIDATION METRICS

- Master Table achieved over 90% completeness across critical fields.
- High consistency was maintained for date, text, and numeric formats.
- Referential integrity across learner and cohort datasets exceeded 85% match rates.

Validation Metric	Target Threshold	Achieved	Remarks
Critical Fields Completeness	≥ 90%	✓ 91.5%	Good coverage; minimal missing values
Text Standardization (Casing, Trimming)	100%	✓ 100%	All text fields normalized
Date Format Consistency	≥ 95%	✓ 97%	Validated with regex checks
Duplicate Record Removal	100%	✓ 100%	All duplicates removed using DISTINCT()
Referential Integrity (Joins across Tables)	≥ 85% match rate	✓ 88%	Learner–Opportunity–Cohort joins successful
Marketing Data Attribution	NA	✗ Not Available	No learner IDs in marketing dataset
Outlier Handling (Demographics)	Identified/Flagged	⚡ Flagged	Outliers flagged, preserved for authenticity

CHALLENGES VS RESOLUTIONS

CHALLENGES VS RESOLUTIONS

Challenge	Resolution
 Missing data	 NULLs preserved; critical fields flag action
 Duplicate entries	 Cleaned using DISTINCT validation methods
 Marketing disconnection	 Separate dashboards with future integration
 Text inconsistency	 Casing normalized across datasets



IMPACT SUMMARY

IMPACT SUMMARY

BEFORE VS AFTER DATA TRANSFORMATION

BEFORE



Fragmented and isolated datasets



33%+ missing demographic data



Duplicate and inconsistent entries



Poor marketing attribution

Manual and error-prone ETL processes

AFTER



Integrated Master Table covering 5 datasets



Over 90% completeness achieved



Data standardized, duplicates removed



Actionable insights for marketing and engagement

Automated, scalable ETL ready for real-time analytics

NEXT STEPS & FINAL TAKEAWAYS

Next Steps -

- Implement continuous Data Quality Monitoring pipelines.
- Strengthen marketing attribution models through improved tracking.
- Extend ETL automation to near real-time processing.
- Leverage cleaned datasets for Predictive Analytics and Machine Learning models.

Final Takeaways -

- The data environment is now trustworthy, consolidated, and strategically positioned for business intelligence.
- The ETL pipeline ensures scalability, flexibility, and sustainability.
- Key systemic data issues have been identified, addressed, and documented.
- Recommendations for marketing-data improvements are critical for full customer journey visibility.



THANK
YOU