

Title: A Modular Framework for Data Harmonization: Enhancing Quality and Efficiency in Healthcare ETL Pipelines

PRESENTER: Freija Descamps

BACKGROUND:

- Harmonizing healthcare data is key for research and interoperability, but the development process is often slowed down by heterogeneous source data, project-specific setup needs and data quality issues.

METHODS

- Our framework automates ETL development by generating project templates based on user input. It supports multiple databases (e.g., PostgreSQL, SQL Server) and input types (e.g., CSV, other databases), and allows users to choose their preferred development approach (Pandas, SQLAlchemy, or raw SQL).

It includes:

- Pre- and post-processing modules for standardizing formats (e.g., dates, nulls)
- Automated scaffolding for customized, consistent project code
- Dockerized local environments with source, target, and lookup schemas
- Ready-to-use local test setup with a dedicated Docker database and preconfigured test code
- Auto-generated data classes for all tables
- Comprehensive logging and metrics to monitor data quality
- Issue tables after each ETL run to detect person-level inconsistencies while preserving privacy

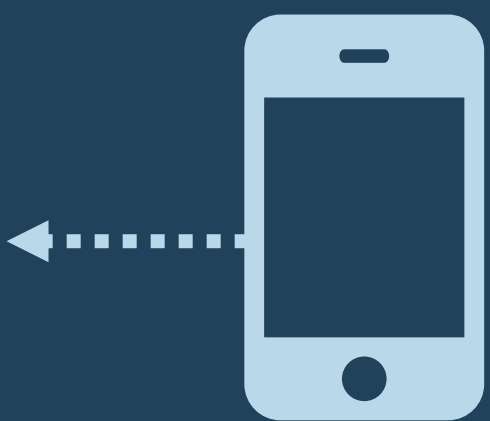
A Modular Framework That Reduces ETL Development Time and Boosts Data Quality in Healthcare Pipelines

HOW IT WORKS:



CONCLUSION

Our framework streamlines healthcare data harmonization by automating key ETL tasks. It reduces development time, improves data quality, and ensures consistency—making ETL pipelines easier to build, maintain, and share across teams.



Take a picture to download the full paper

Isaac Claessen¹, Silvia Jimenez Navarro¹, Shirah Cashriel¹, Panagiotis Gialernios¹
¹edenceHealth NV

