leap**logic**

an **IMPETUS** product

# Automated transformation of DataStage ETL to cloud-native equivalent or PySpark

Modernizing legacy ETL scripts to modern cloud-native platforms or open collaboration-based languages has become a strategic imperative for enterprises struggling with petabytes of unstructured and fast data from multiple sources. Enterprises are looking to move from DataStage because of the high cost of ownership, limited documentation and cloud integration capabilities, and complex set up process.

However, migrating DataStage ETL to modern data architecture is complex as it is difficult to edit columns and mapping between the stages. The unavailability of source code versioning also makes it difficult to roll back to previous versions in case of any errors, where logs are not handy and difficult to decipher and debug. Businesses, therefore, have several concerns:

- Will there be any business downtime?
- How do I transform years of complex business logic and code?
- Will my workloads be optimized to address the nuances of the new environment?
- How do I prioritize the workloads for transformation?
- How do I ensure seamless operationalization of ETL on the target environment?
- How do I validate the transformed code and migrated data?

## Key Benefits

Compared to traditional approaches, LeapLogic enables:

- **4x** faster transformation
- **2x** lower cost
- **4x** developer productivity
- **100%** cloud-ready
- **100%** automation across the migration lifecycle

LeapLogic, an Impetus product for automated workload transformation addresses all these concerns. Its intelligent grammar engine identifies optimization opportunities at schema, code, and execution level and automatically converts all types of workloads, logic, and workflows to Snowflake-native equivalent, Databricks-native equivalent, or PySpark.

## Key Features

- Intelligent assessment and recommendation for the target architecture and tech stack

- Business logic conversion

- End-to-end transformation to PySpark or a cloud-native equivalent

- End-to-end packaging, orchestration, and execution for the target

- Code optimization and query validation to avoid business disruption

- Automated legacy code translation to PySpark with multiple query engine support

- Cost-performance ratio optimization

- Data governance and security compliance

# How it works

LeapLogic enables end-to-end transformation, operationalization, and transitioning of workloads in four steps:

LEGACY DATA PLATFORM

IBM DataStage

**ASSESS**
- Inventory listing
- Lineage–dependancy analysis
- Capacity planning
- Resource estimation

**TRANSFORM**
- Auto-transform DataStage job components and functions, data
- Packaged using cloud-native wrappers/open collaboration languages such as PySpark code
- Repeatability, extensibility

LOGIC TRANSFORMATION

**VALIDATE**
- Pipeline-based validation
  - Schema, metadata, data
  - Data-based code
  - Pre and post processed data

**EXECUTE**
- Executable package with cloud-native orchestrators/ PySpark code
- CI/CD and transition support

MODERN DATA PLATFORM
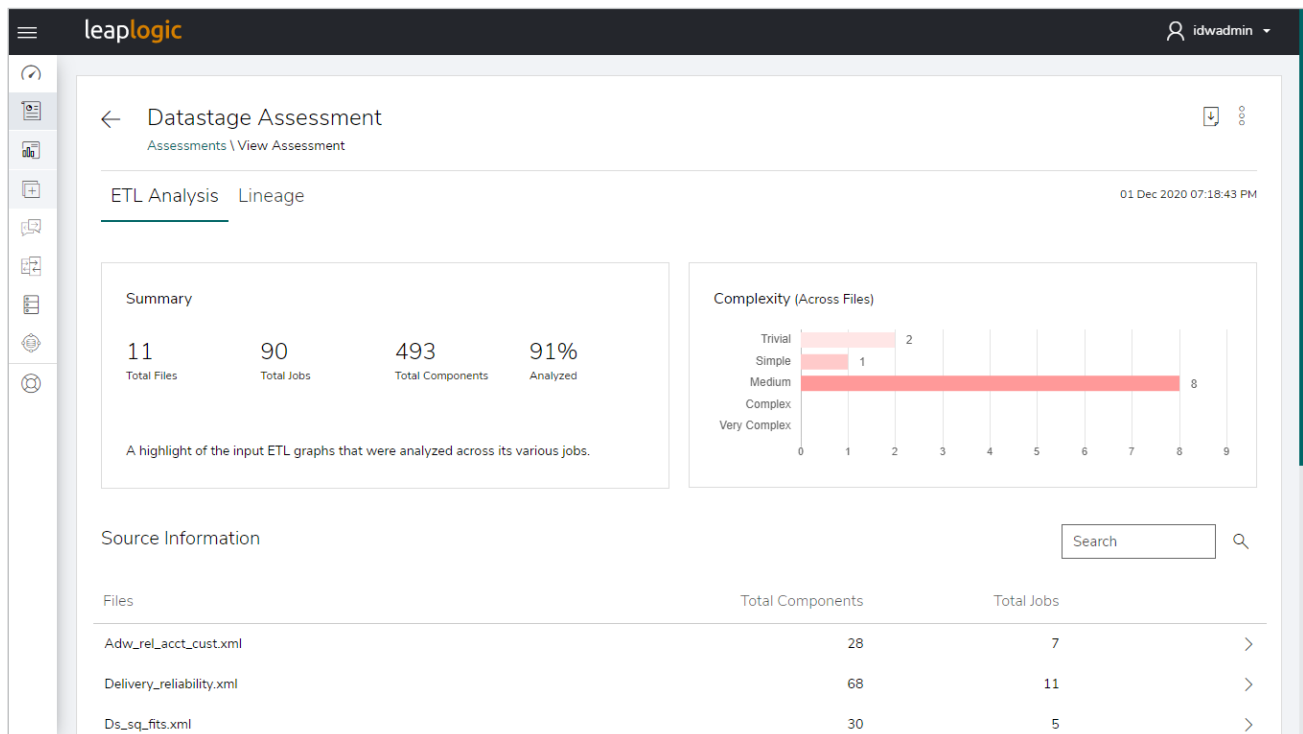
databricks     PySpark

Spark     snowflake

Automated workload transformation

# Step 1: Assessment and prescription

- Lists entire inventory for diverse workloads

  - Assesses ETL scripts, DML and DDL scripts, procedures, scheduler/orchestrator scripts (jobs), etc.

  - Identifies total jobs and components per ETL script

  - Identifies job activities, parallel jobs, sequential jobs, aggregator functions, joins, filters, execution command activities, exception handler, etc.

  - Provides actionable insights and prescriptive recommendations with detailed downloadable reports

- Identifies complex interdependencies to group workloads for offload

  - Identifies job-level complexity and its aggregate across files

  - Plots comprehensive, graphical lineage of DataStage XMLs across jobs highlighting the participating input and output tables

  - Provides advanced filters according to workload type and an interactive graphical interface to deep dive into certain flows



DataStage Jobs Assessment



DataStage Jobs Assessment

# Step 2: Transformation

- Transforms DataStage ETL scripts and migrates schema and dataset to the target store of choice

- Native conversion, packaging, and orchestration

  - Transforms core business logic to cloud-native wrappers or orchestrators

  - Re-packages business logic with scripts to production-ready jobs

  - Generates end-to-end executable package for on staging and production environments after a thorough system integration testing

- Extensible tool and methodology

  - Extensible, repeatable, and verifiable methodology

  - Converts code to a variety of target stores and formats, enabling enterprise-wide usage



Converted PySpark code



Transformation of DataStage components

## Step 3: Validation

- Pipeline-based validation of transformed code with little manual intervention
  - Instantly verifies the transformed jobs and components

## Step 4: Operationalization

- Delivers targret-specific executable package
  - Cloud-native orchestration and execution on production
- Supports operationalization
  - Supports end-to-end transitioning into production and operationalization
  - Optimizes target capacity
  - Stabilizes environment through a parallel-run period
  - Smooth cut-over planning
  - Ensures implicit data governance and compliance on the cloud
  - Ensures continuous integration and delivery
  - Helps in enabling powerful operational monitoring on target
  - Provides runbook documentation, training, and handholding

To start your automated transformation from DataStage to cloud-native equivalent or PySpark write to us at: **info@leaplogic.io**

BOOK A DEMO ↗

## leaplogic