# Pós-graduação em Data Science & Business Analytics Formato Blended 2ª Edição - 2022

ISEG Executive Education

**Data Warehousing**

**Artur Vieira**

# Índice

- Traditional Data Warehouse Architectures

- Data Integration

- SQL Server Integration Services

- Azure Data Factory

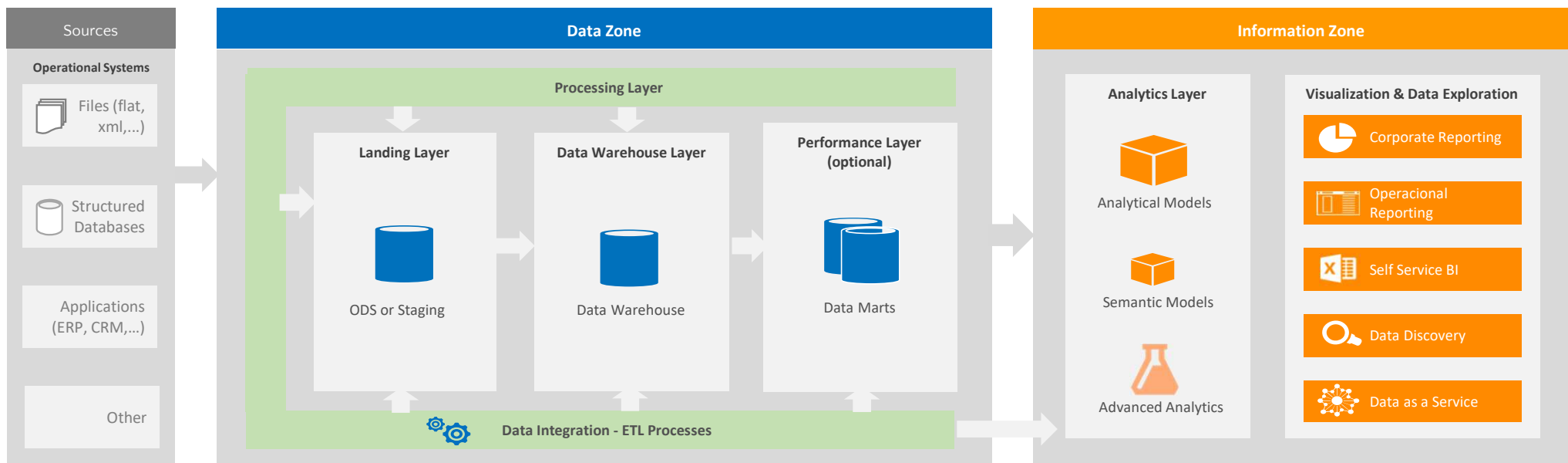# Traditional Data Warehouse Architecture

| Sources | Data Zone | | Information Zone |
|---|---|---|---|

**Operational Systems**

Files (flat, xml,...)

Structured Databases

Applications (ERP, CRM,...)

Other

**Processing Layer**

**Landing Layer**

ODS or Staging

**Data Warehouse Layer**

Data Warehouse

**Performance Layer (optional)**

Data Marts

**Data Integration - ETL Processes**

**Analytics Layer**

Analytical Models

Semantic Models

Advanced Analytics

**Visualization & Data Exploration**

Corporate Reporting

Operacional Reporting

Self Service BI

Data Discovery

Data as a Service

ODS – operational data store

# Traditional Data Warehouse Architecture

ODS – operational data store
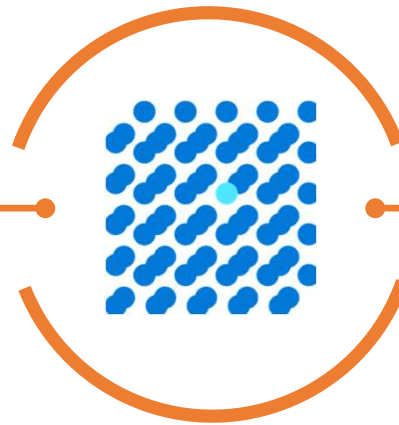
# Tópico    Data Integration

**Data integration is a process in which heterogeneous data is retrieved and combined as an incorporated form and structure.**

- **Extract, Transform and load (ETL)**

- Integrate structured and unstructured data

- Multiple sources

- Multiple destinations

- Data Modeling

- Data profiling

- Data Cleansing, Data Merging / Data Enrichment

# Access to data remains top issue

**Less than half of structured data** is actively used in decision-making

**Less than 1% of the unstructured data** is analyzed or used

**97%\* of executives** find data silos harmful to their organization

*\*83% of executives confirm their organizations have data silos*

Harvard Business Review, 2017:
https://hbr.org/2017/05/whats-your-data-strategy

American Management Association
2017 survey

# There are barriers to getting value from data

Tópico

⚠️ Data silos

⚠️ Incongruent data types

⚠️ Complexity of solutions

⚠️ Multi cloud environment

⚠️ Rising costs

Tópico

# Derive real value from your data

| Data silos | Incongruent data types | Complexity of solutions | Multi cloud environment | Rising costs |
|---|---|---|---|---|

| One hub for all data | Support for diverse types of data | Unlimited data scale | Familiar tools and ecosystem | Lower TCO |
|---|---|---|---|---|

On-premises, hybrid, Azure

# SQL Server Integration Services (SSIS)

SQL Server Integration Services (SSIS) is a component of the Microsoft SQL Server database software that can be used to perform a broad range of data integration and data transformation task.

- **Data integration -** it combines the data residing in different sources and provide users with a unified view of these data

- **Data transformation - it transforms the ingested data by  applying logic to fit the data objectives**

# Azure Data Factory

### Connect with confidence

All-inclusive connectivity that prioritizes security and compliance

### Reduce integration costs

Serverless, scales on demand to focus on the data, not infrastructure

### Work efficiently

Intuitive, visual environment for everyone

## Productive & trusted hybrid data integration service that simplifies ETL with any data, from any source, at scale.

Tópico

# Connect with confidence

## All-inclusive connectivity

More than 80 natively built and fully managed connectors, no added cost, new connectors added monthly

Efficient and resilient data transfer by leveraging the full capacity of underlying network bandwidth, up to 2 GB/sec throughput

## Trusted, global cloud presence

Data Factory availability in 25+ regions, with data movement available globally to help ensure compliance & reduced network egress costs.

## Security & compliance peace of mind

Native integration with Azure Active Directory (AAD) and Azure Key Vault (AKV) for identity and access management to cloud solutions & applications, based on centralized policy and rules

HIPAA, HITECH, ISO/IEC 27001, ISO/IEC 27018, CSA STAR certification.

New Dataset

| | | | |
|---|---|---|---|
| HubSpot | Google Big Query | Jira Software | Magento |
| Marketo | eloqua | amazon web services | Adobe Analytics |
| Acumatica The Cloud ERP | amazon REDSHIFT | Azure Audit Logs | Azure Mobile Engagement |
| cloudera IMPALA | GitHub | Visual Studio | SendGrid |
| webtrends | amazon S3 | salesforce | Google Analytics |

# Reduce integration costs

**Serverless, fully managed service**

No infrastructure to manage, no hardware to upgrade
Scales on demand
Pay only for what you use.

**One data integration service for everyone**

Reduce integration tool fragmentation & costs
Flexibility to work how you please, visually or using code
(Python, .NET or ARM)

**Fast and scalable transformations with Spark**

Azure Databricks' Spark engine powers data
transformations for fast and fully managed data
transformations
.

**Reduce development overhead**

Migrate to the cloud by moving SSIS packages into Azure
without redevelopment
Use existing tools for new development.
Full integration with GitHub for team collaboration.

Azure

SaaS

Enterprise
apps and data

Clouds

# Work efficiently

**Simple to get started**

**Azure Data Factory Dashboard**: Use tutorials, quick starts, predefined <u>templates</u>, leverage & share best practices & patterns.

**Easy to be productive**

**Visual environment**: Ingest, move, prepare, transform and process your data with just a few clicks.

*   **Data orchestration**: Visually <u>construct workflows</u> to orchestrate integration and transformation.

*   **Data transformation**: <u>Mapping Data Flows</u> to visually create complex pipelines and transforms. Native handling of data evolution / schema drift & for non–relational data, Rich & granular monitoring and management

*   **Pipeline automation**: Automate pipeline runs with <u>Triggers</u>

*   **Intelligent Data preparation**: Visually explore data with <u>Wrangling Data Flows</u>

*   **CI/CD**: Simple dev ops integration with built in support with Azure Monitor, API, PowerShell, Azure Monitor logs, and health panes on the Azure Portal, Git integration

# Azure Data Factory - Data Integration Service

Azure Data Factory (ADF) is a cloud-based data integration service **that orchestrates and automates the movement and transformation of data**.

**It orchestrates existing services** that collect raw data and transform it into ready-to-use information. ADF is used to **collect data from many different data sources, ingest and prepare it, organize and analyze it with a range of transformations, then publish ready-to-use data for consumption.**

# Azure Data Factory – Control Flow

Coordinate pipeline activities into **finite execution steps to enable looping, conditionals and chaining while separating data transformations into individual data flows**

# Azure Data Factory - Data Integration Service



## Data Factory

A data integration account.
Location of orchestration, service metadata

## Integration Runtime (IR)

ADF's execution engine

Three core capabilities:
- data movement
- pipeline activity execution
- SSIS package execution

# What are Mapping Data Flows?

**Data Flow is a new feature of Azure Data Factory to build data transformations in a visual user interface**

- Transform at scale, in the cloud
- Code-free pipelines do NOT require understanding of Spark / Scala / Python / Java
- Serverless scale-out transformation execution engine
- Resilient data transformation Flows built for big data scenarios with unstructured data requirements
- Operationalized with Data Factory scheduling, control flow and monitoring

# Mapping & Wrangling Data Flows



**MAPPING DATAFLOW**
Code-free data transformation @scale

**WRANGLING DATAFLOW**
Code-free data preparation @scale

PUBLIC PREVIEW

# Modern Data Warehouse (MDW)

**On-premises data**

Oracle, SQL,, Teradata, fileshares, SAP

**Cloud data**

Azure, AWS, GCP

**SaaS data**

Salesforce, Dynamics

**INGEST**

Azure
Data Factory

**PREPARE**

Azure
Data Factory

**Azure
Databricks**

**TRANSFORM,
PREDICT
& ENRICH**

Azure
Data Factory

**Azure
Databricks** **Azure ML**

**SERVE**

Azure
SQL Data
Warehouse

**VISUALIZE**

Power BI

**STORE**

**Azure Data Lake Storage Gen2**

**Data Pipeline Orchestration & Monitoring**

**Azure Data Factory**

# Modern Data Warehouse Pattern Today



Databases

Logs, files, and media (unstructured)

Business/custom apps (structured)

**Data Loading**

Azure Data Factory

**Ingest storage**

Azure Storage/ Data Lake Store

Load flat files into data lake on a schedule

**Data processing**

Azure Databricks

Read data from files using DBFS

**Serving storage**

Azure SQL DW

Load processed data into tables optimized for analytics

Applications

Power BI

Dashboards

**Orchestration**

Azure Data Factory

Extract and transform relational data

Clean and join with stored data

Load to SQL DW

# Modern Data Warehouse Pattern with Mapping Data Flows



**Data Loading**

Azure Data Factory

**Ingest storage**

Azure Storage/ Data Lake Store

Load files into data lake on a schedule

**Data Flow Data Transformation**

Azure Data Factory

Azure Databricks

Extract and transform relational data

Clean and join disparate data

**Serving storage**

Azure SQL DW

Load processed data into tables optimized for analytics

Databases

Logs, files, and media (unstructured)

Business/custom apps (structured)

Applications

Dashboards

Power BI

**Scheduled & orchestrated by ADF**

# Pipeline execution of a Data Flow Activity



- Design code-free ETL workflows
- Copy data from on-prem, other clouds and Azure
- Stage data for transformation
- Build visual data transformations
- Schedule triggers for your pipeline execution
- Monitor processes and configure alerts
- All within ADF

# Best in class monitoring and management

Monitor Pipeline and Activity Runs

Rich language to query Runs

Operational lineage between parent-child pipelines

Azure Monitor Integration

- Diagnostics logging
- Metrics & Alerts
- Events

Restate Pipeline and Activities

# Data Transformations, Expression Language and Debugging

# Azure Data Factory Continues to Extend Data Flow Library

# Expression builder

# Switch to Debug Mode and select sample data to work with for debugging

- Set Parameter values and sample data in debug settings
  - Change # of rows used per source
  - Replace source with debug dataset
  - Assign debug parameter values

# Debug Data Flows with Data Preview and Data Sampling

# Mapping Data Flow common scenarios

# Slowly Changing Dimension Scenario



- Common DW pattern to manage changing attributes to dimension members
- Graphically build code-free SCD ETL pattern to load your data warehouse
- Connect directly to Azure SQL DB and Azure SQL DW
- Use Lookup, Surrogate Key, Derived Column and Select transforms

# Data De-Duplication



- Use this pattern to eliminate common rows from your data
- You pick a heuristic to use during duplicate matching
- You can tag rows and/or remove duplicate rows
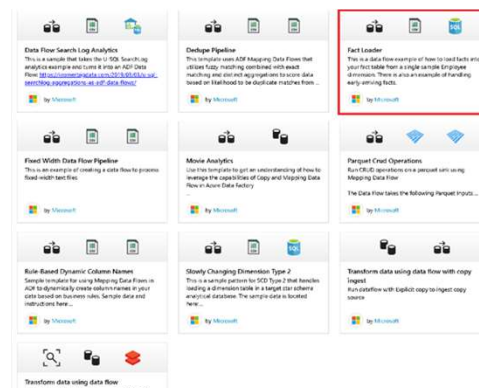- Use exact matching and/or fuzzy matching
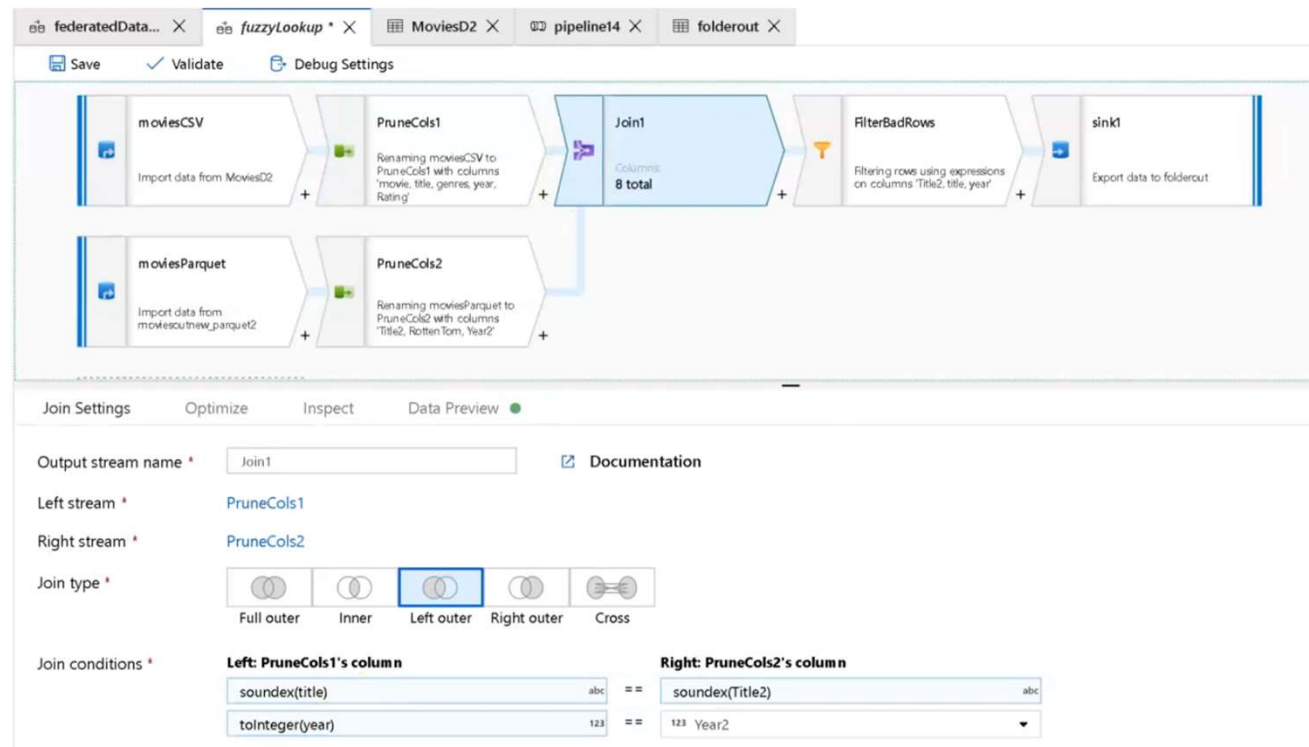- Available as pipeline template *Dedupe Pipeline*

# Load Fact Table in DW Scenario



- Classic ETL pattern is easy to build in ADF's code-free Data Flow visual data transformation environment

- Add Aggregate transforms to produce calculations that you store in your analytical database schema

- Use Join transform to combine data from multiple data sources and data streams inside your data flow

- Land your data in your Lake folders or direct to Azure SQL DW

# Fuzzy Lookups

- Sometime when performing inline lookups, you don't have exact matches when looking for references
- Fuzzy Lookups with Soundex helps find matches based on phonetic algorithms
- Very useful in data lake scenarios where joins and lookups are against data that is not normalized or cleaned

# Data Lake Data Science Scenario



- ADF supports building visual data transformations against your data directly in Data Lake locations (i.e. Azure Blob Store, Azure Data Lake Store)

- Built-in handling of schema drift for frequent changes in data lake file formats, columns, and data types

- Perform data exploration and data profiling across your data lake in ADF Data Flow with interactive debug data preview and quick actions

# Schema Drift

# Schema Drift

In most real-world data integration solutions, source and target data stores will change shape

- Source data fields will change name
- Number of columns will change over time

Traditional ETL processes break when schemas drift
Mapping Data Flow has built-in facilities for flexible schemas to handle schema drift

- Patterns, rule-based mapping, byName function
- Source: Read additional columns on top of what is defined in the dataset source
- Sink: Write additional columns on top of what is defined in the dataset sink

**Pattern matching**

· Match by name, type, stream, ordinal position

## Rule-based mapping

· Rather than pick and choose columns for transformations one–by–one, build policies that collect columns based on matching rules.

# Prepare Practical Classes

# Prepare Practical Classes

- Access the URL
https://portal.azure.com/#@isegulisboacloud.onmicrosoft.com/resource/subscriptions/000a9103-fec4-4ce8-aba5-d7b0908f5aca/resourceGroups/DSBAFB2_DW/overview

- User your Power BI student Account ex. ixxxxx@students.isegexecutive.education
- Accept the access to Azure Portal resources
- Validate you are able to see the resource group DW

# Resources

- Book: UNDERSTANDING AZURE DATA FACTORY, Rawat, Sudhir, Narain, Abhishek
- Patterns: http://aka.ms/dataflowpatterns
- Documentation: https://docs.microsoft.com/en-us/azure/data-factory/concepts-data-flow-overview

**Executive Education**