

Introducción a IBM DataStage

Panorama del proceso ETL y arquitectura

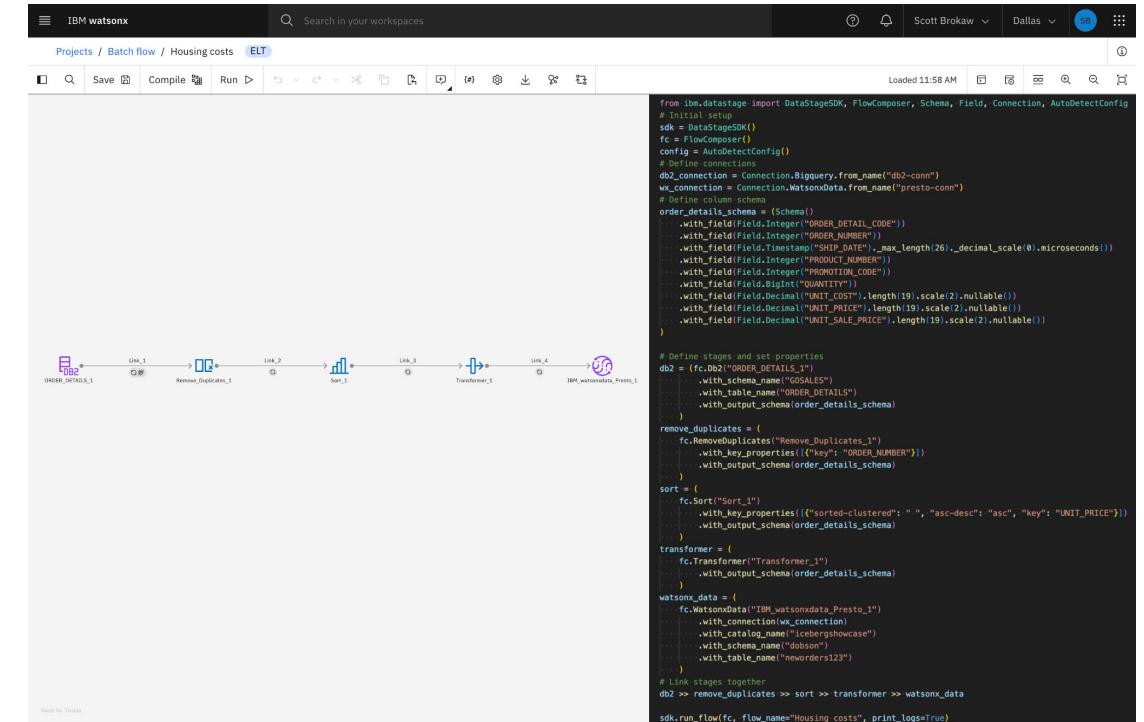
Expositor: Consultor Senior BI

Fecha: Febrero 2024

¿Qué es IBM DataStage?

IBM InfoSphere DataStage es una herramienta líder en integración de datos para:

-  Diseñar procesos ETL complejos.
-  Integrar múltiples fuentes heterogéneas.
-  Cargar Data Warehouses y Data Lakes.
-  Soporta procesamiento paralelo masivo.



The screenshot shows the IBM WatsonX DataStage interface. At the top, there's a navigation bar with 'IBM watsonx', 'Search in your workspaces', and user information like 'Scott Brokaw' and 'Dallas'. Below the bar is a toolbar with icons for 'Save', 'Compile', 'Run', and others. The main area has tabs for 'Projects', 'Batch flow', 'Housing costs', and 'ELT'. A large code editor window displays Python code for an ETL flow. The code imports DataStageSDK, FlowComposer, Schema, Field, Connection, AutoDetectConfig, and defines connections for db2 and wx. It then sets up stages: 'Remove_Duplicates_1' (using 'fc.RemoveDuplicates("Remove_Duplicates_1")'), 'Sort_1' (using 'fc.Sort("Sort_1")'), and 'Transformer_1' (using 'fc.Transformer("Transformer_1")'). Finally, it links stages together and runs the flow.

```
from ibm_datastage import DataStageSDK, FlowComposer, Schema, Field, Connection, AutoDetectConfig
# Initial setup
sdk = DataStageSDK()
fc = FlowComposer()
config = AutoDetectConfig()
# Define connections
db2_connection = Connection.Bigquery.from_name("db2-conn")
wx_connection = Connection.WatsonData.from_name("presto-conn")
# Define schema
order_details_schema = (Schema()
    .with_field(Field.Integer("ORDER_DETAIL_CODE"))
    .with_field(Field.Integer("ORDER_NUMBER"))
    .with_field(Field.Timestamp("SHIP_DATE"), max_length(26), decimal_scale(0), microseconds())
    .with_field(Field.Integer("PRODUCT_NUMBER"))
    .with_field(Field.Integer("UNIT_COST"))
    .with_field(Field.BigInt("QUANTITY"))
    .with_field(Field.Decimal("UNIT_COST", length(19), scale(0), nullable()))
    .with_field(Field.Decimal("UNIT_PRICE", length(19), scale(2), nullable()))
    .with_field(Field.Decimal("UNIT_SALE_PRICE", length(19), scale(2), nullable())))
)
# Define stages and set properties
db2 = (fc.db2("ORDER_DETAILS_1")
    .with_schema_name("GOSALES")
    .with_table_name("ORDER_DETAILS")
    .with_output_schema(order_details_schema)
)
remove_duplicates = (
    fc.RemoveDuplicates("Remove_Duplicates_1")
    .with_key_properties([("key", "ORDER_NUMBER")])
    .with_output_schema(order_details_schema)
)
sort = (
    fc.Sort("Sort_1")
    .with_key_properties([("sorted-clustered", "", "asc-desc", "asc", "key": "UNIT_PRICE")])
    .with_output_schema(order_details_schema)
)
transformer = (
    fc.Transformer("Transformer_1")
    .with_output_schema(order_details_schema)
)
watsonx_data = (
    fc.WatsonData("IBM_watsondata_Presto_1")
    .with_connection(wx.connection)
    .with_catalog_name("icebergshowcase")
    .with_schema_name("dobsom")
    .with_table_name("neworders123")
)
# Link stages together
db2 >> remove_duplicates >> sort >> transformer >> watsonx_data
sdk.run_flow(fc, flow_name="Housing costs", print_logs=True)
```

¿Para qué se utiliza?



Integración Empresarial

Consolidación de datos de toda la organización en un punto único de verdad.



Migración de Datos

Movimiento seguro de información entre sistemas legacy y plataformas modernas.



Business Intelligence

Construcción de bases sólidas para reporting, analítica avanzada y Big Data.

Concepto Fundamental: ETL

E Extract (Extraer)



T Transform (Transformar)



L Load (Cargar)



Proceso clave para convertir **datos crudos** en **información confiable** para la toma de decisiones empresariales.

Fases del Proceso ETL

1. Extracción

- Bases de Datos (Oracle, DB2)
- Archivos (CSV, JSON)
- APIs y CRM/ERP

2. Transformación

- Limpieza y Calidad
- Reglas de Negocio
- Uniones y Cálculos

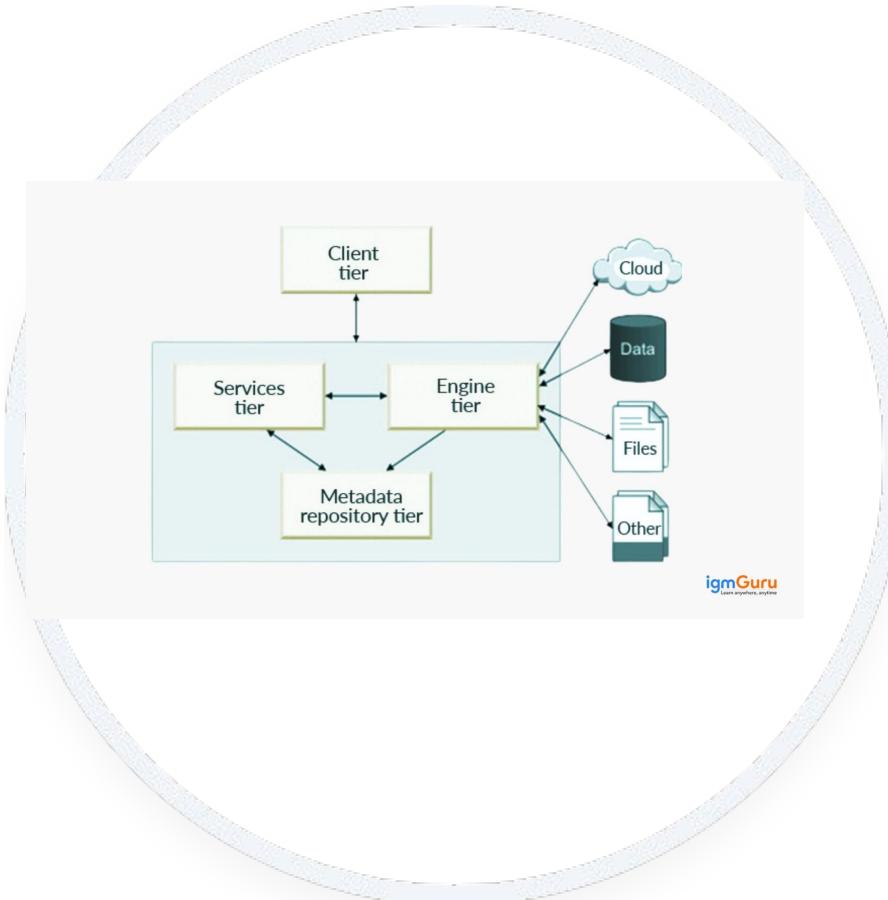
3. Carga

- Data Warehouses
- Data Marts
- Sistemas Destino

Flujo Lógico del Proceso



El Rol de DataStage



- 💡 **Diseño Gráfico:** Interfaz visual para procesos complejos.
- ⚙️ **Procesamiento Paralelo:** Optimización de recursos.
- 🕒 **Automatización:** Programación de cargas recurrentes.
- 🔍 **Monitoreo:** Control total de la ejecución.



Arquitectura

Componentes, Motor y Repositorio

Vista Lógica de la Arquitectura

El flujo de información en la infraestructura:



Componentes Principales



1. Capa Cliente

Interfaz gráfica para diseño de Jobs, configuración de parámetros y etapas.



2. Capa Servidor

Administración de la ejecución, seguridad y gestión de conexiones concurrentes.



3. Engine (Motor)

Núcleo que ejecuta las transformaciones con procesamiento paralelo optimizado.

Tipos de Jobs

Parallel Jobs

Diseñados para alto rendimiento. Dividen la carga en múltiples nodos de procesamiento.

Server Jobs

Arquitectura secuencial más antigua. Útiles para procesos ligeros o legados.

Sequence Jobs

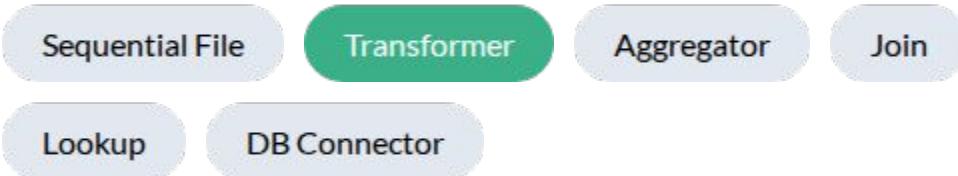
Orquestadores que controlan el flujo y las dependencias entre múltiples jobs.

Unidades de Trabajo

-  **Jobs:** La unidad básica de ejecución ETL.
-  **Links:** Conexiones que transportan los datos entre Stages.
-  **Director:** Herramienta de monitoreo y logs.
-  **Stages:** Bloques de función específica (Join, Sort, etc).
-  **Repository:** Almacén central de metadatos y versiones.

¿Qué es un Stage?

Componentes modulares dentro de un Job que ejecutan acciones sobre los datos:



Stage: Transformer

Es el corazón de la lógica de negocio en DataStage.

Permite:

- ✓ Cálculos matemáticos y lógicos.
- ✓ Aplicación de condiciones IF/THEN.
- ✓ Creación de columnas derivadas.
- ✓ Conversión de tipos de datos complejos.



Combinación de Datos

Join Stage

Combina múltiples flujos de grandes volúmenes.

Equivale al JOIN de SQL. Requiere datos ordenados.

Lookup Stage

Ideal para validar contra tablas de referencia en memoria. Optimizado para "Left Outer Joins".

Metadatos: Repository

El repositorio central almacena:

-  Definiciones de Tablas.
-  Diseños de Jobs y Rutinas.
-  Control de Versiones.
-  Parámetros Compartidos.



Operación: Director



Herramienta fundamental para el equipo de Operaciones:

- ▶ Ejecución manual o programada.
- ✖ Análisis profundo de logs.
- ⚠ Diagnóstico y resolución de errores.
- ⟳ Reinicio selectivo de procesos.

Procesamiento Paralelo



DataStage utiliza **Particionamiento de Datos** para ejecutar tareas simultáneas, garantizando escalabilidad lineal en entornos de Big Data.

¿Por qué elegir DataStage?

- ✓ **Escalabilidad:** Crece con tus volúmenes de datos.
- ✓ **Conectividad:** Soporte nativo para nubes y Big Data.
- ✓ **Estabilidad:** Estándar de la industria financiera.
- ✓ **Gobernanza:** Control centralizado de metadatos.



**Líder en el
Cuadrante de Gartner**

Reconocido mundialmente por su robustez.

Conclusiones

- ✓ Herramienta ETL empresarial robusta.
 - ✓ Arquitectura escalable Cliente-Servidor-Motor.
 - ✓ Dominio del paralelismo para grandes volúmenes.
 - ✓ Basado en Jobs, Stages y Repositorio.
-

¿Preguntas?

Muchas gracias por su atención.

www.grupcief.com