# Credits

**Prepared by**
Data Migration Jumpstart Engineering Team
Advanced Solutions Delivery – Data & AI CTO

askdmjfordmtools@microsoft.com
dmjarchitects@microsoft.com

Revision 2
3/26/2020

# Disclaimer

Data Migration Jumpstart Engineering
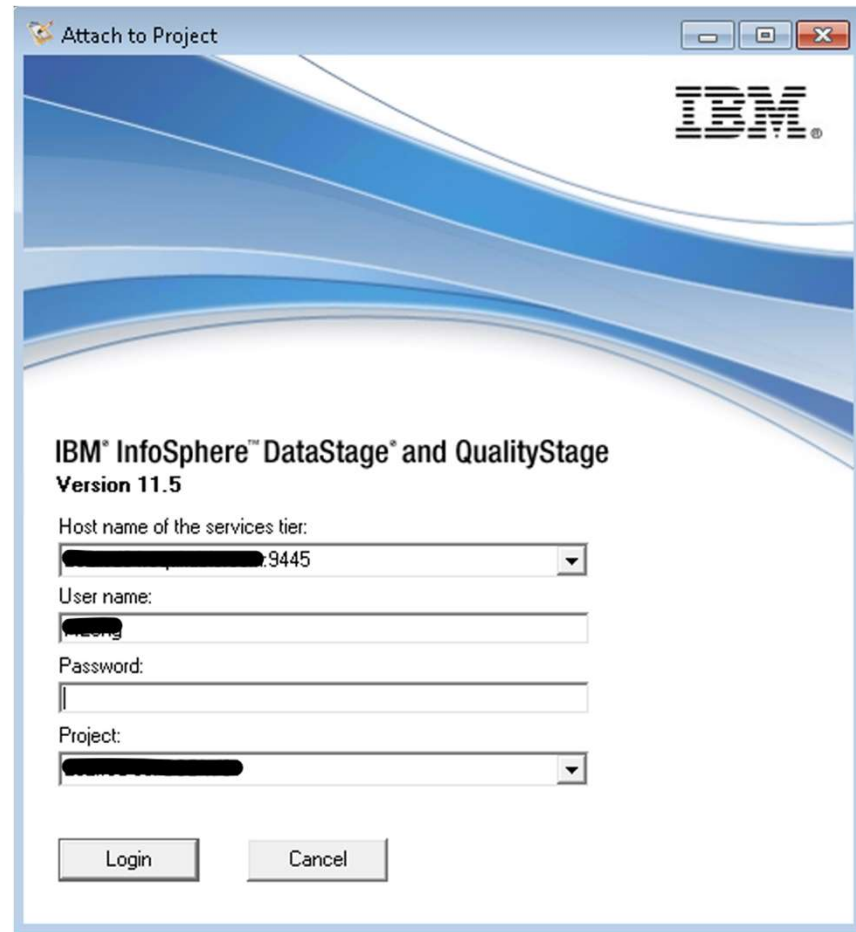
# Topics

- [Utilizing DataStage](#)
- [Singleton and Bulk Load Inserts](#)
- [Insert and Update data within Azure Synapse Analytics with DataStage](#)
- [Changes to Transformers](#)
- [Executing Azure Synapse Analytics Stored Procedures with DataStage](#)
- [Large Data Ingestion in Azure Synapse Strategy](#)
- [COPY INTO statement sample](#)

Data Migration Jumpstart Engineering

# Utilizing DataStage

- A client DataStage application is installed on workstation/laptop (Designer)
- Clicking on the Designer icon opens the DataStage application which requires authentication to the server with name, password and project name
- Selection of a job is made from a tree structure on the left showing the jobs and sequences in hierarchical fashion

Data Migration Jumpstart Engineering

# IBM InfoSphere DataStage Logon Screen

# DataStage Initial Screen

# Select Job from the Project

# Working with the DataStage palette

- On the left of the screen there are selections for the *palette*.
- The *Database* category allows the user to select the *connectivity* type for this job.
- The *Processing* category allows the user to select the type of processing to occur, i.e. *transformer*
- The *File* category allows the user to select a type of target or source file (as opposed to a database) and requires no connectivity information

# Sample Netezza Job – Update and Insert

# Selection for Connectivity of Available Sources/Targets

# Singleton and Bulk Load Inserts

- Selection of **ODBC** Connector will only allow singleton transactions
- Selection of **DRS** Connector will allow for selection of Bulk Insert option
- All stored procedures run on Azure Synapse Analytics must be executed via the Stored Procedure Database option from this selection screen

# DRS Connection

- Use the DRS Connector stage to access relational database management systems by using the native interfaces that are available for the corresponding databases.
- The choice of the database type is not determined when the stage is placed on the job canvas but is instead specified when the stage is configured.
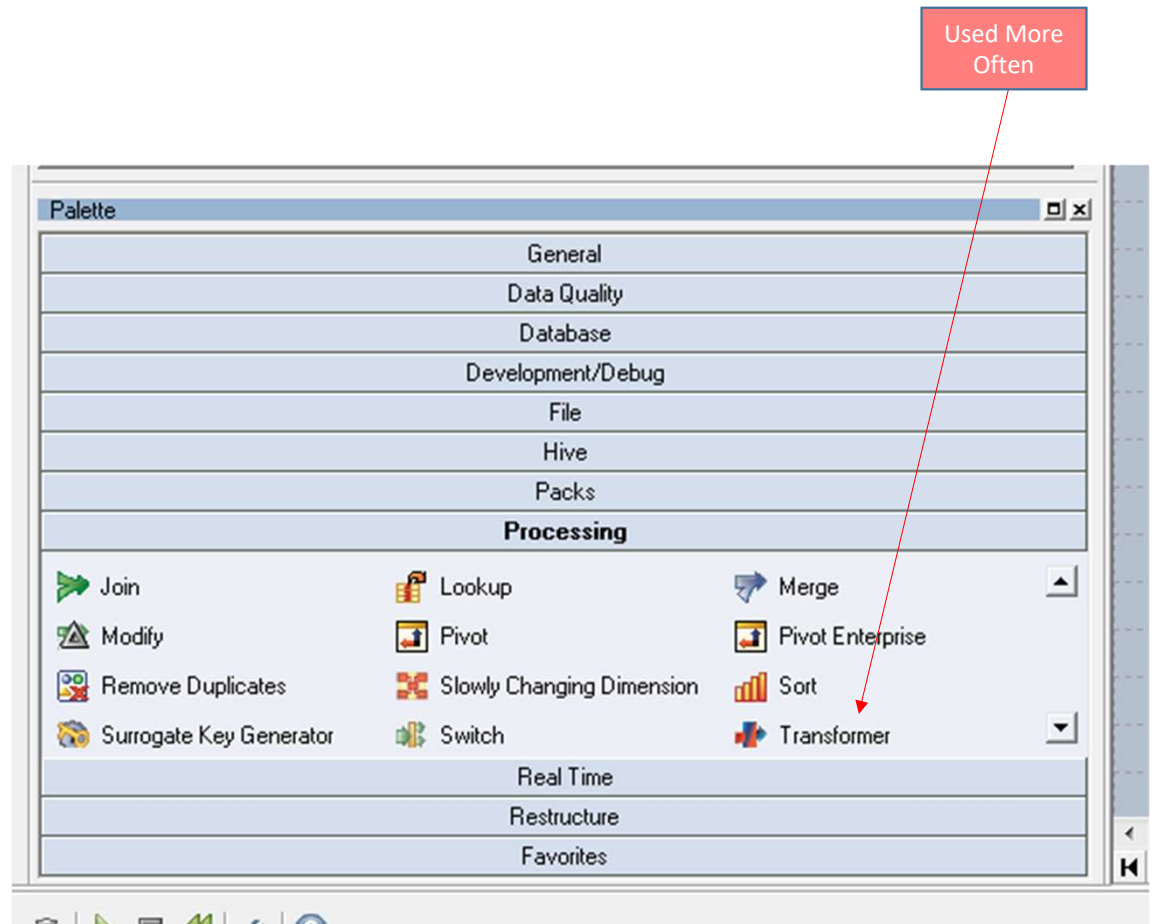- The DRS Connector stage supports IBM® DB2®, Oracle, and ODBC data sources. For other database types, you can configure the DRS Connector stage to use the ODBC database type and access the databases through the ODBC drivers that are included with InfoSphere Information Server.
- After the DRS Connector stage is placed on the job canvas, it needs to be configured to perform the operation intended by the job design.
- When a new DRS Connector stage is added to the server or parallel canvas, you must specify the connection information for the database that the stage connects to at run time.

- Information Links
    - DRS Connection Overview
    - DRS Configuration
    - Settings for the DRS Connector stage
    - Defining a DRS Connector stage connection to the database

# Selection of Choices for Processing



Used More Often

Data Migration Jumpstart Engineering

# Selection of File Options

# Setup Database Tables for use

- Any Database table(s) used in this job must be imported into the job
- Every table that is used in the job must be imported individually.
- The name used to set up the ODBC connectivity for the database is utilized in the "DSN Name" box
- Authorization credentials are specified in each individual import stage

# Import Table Selections



Data Migration Jumpstart Engineering

# Import each table used



The database schema must be included in the table import along with the username and password.

# DataStage ODBC.ini driver parameters

- In order to utilize the Bulk Load option, the *EnableBulkLoad* parameter must be set to 1.

```
Driver=/software/IBM/InformationServer/Server/branded_odbc/lib/VMsqls00.so
Description=DataDirect SQL Server Wire Protocol driver
AlternateServers=
AlwaysReportTriggerResults=0
AnsiNPW=1
ApplicationName=
ApplicationUsingThreads=1
AuthenticationMethod=1
BulkBinaryThreshold=32
BulkCharacterThreshold=-1
BulkLoadBatchSize=1024
BulkLoadFieldDelimiter=
BulkLoadOptions=2
BulkLoadRecordDelimiter=
ConnectionReset=0
ConnectionRetryCount=0
ConnectionRetryDelay=3
Database=SQLDatawarehouse_POC
EnableBulkLoad=1
EnableQuotedIdentifiers=1
EncryptionMethod=1
FailoverGranularity=0
FailoverMode=0
FailoverPreconnect=0
FetchTSWTZasTimestamp=0
FetchTWFSasTime=1
GSSClient=native
HostName=sqldatawarehousepoc.database.windows.net
HostNameInCertificate=
InitializationString=
Language=
```

Data Migration Jumpstart Engineering

# Converted Sample Job – Update and Insert

# Utilize same method of pulling information



Data pulled by "user defined SQL" statement.

Data Migration Jumpstart Engineering

# For updating data, first, ingest into staging table using Bulk Insert.



Data Migration Jumpstart Engineering

Detail example of the Update statement.

# Similar approach to Insert



Load into Stage Table

Data Migration Jumpstart Engineering

# Changes to Transformers

- When changing the sources and targets, assure that you are replacing the existing stages by deleting them and then pasting the new ones in.
- Do not delete the connection arrows.
- Open each Transformer processing module and change necessary sources and/or targets information

# Change Transformer Processor

# Executing Stored Proc on Azure Synapse Analytics

- If an additional step is required to creating a sequence
- Create the stored procedure on Azure Synapse Analytics SQL Pool
- Create a job to execute the stored procedure
- Set up a "Sequence" job to execute the job that inserts and updates and then executes the stored proc

# Stored Procedure to Build Key and Insert

```
CREATE PROCEDURE STAGE_INSERT AS
BEGIN
DECLARE @security_id bigint
Select @security_id = MAX(ROW_ID) FROM SAMPLE_SCHEMA.W_BDA_SECURITY_DM_Dist
--
if Exists(select 1 from sys.tables where name = 'W_BDA_SECURITY_DM_4_ROWNUM') Drop Table STAGE.W_BDA_SECURITY_DM_4_ROWNUM
create table STAGE.W_BDA_SECURITY_DM_4_ROWNUM
with (distribution = replicate, heap) as
select (@security_id +row_number() over(order by sec_no, cusip asc)) as NEW_ROW_NUM,
SEC_NO, CUSIP,
SEC_LONG_NAME,
SEC_SHORT_NAME,
SEC_SYMBOL,
SEC_CLASS_CD,
OLD_SEC_TYPE_CD,
NEW_SEC_TYPE_CD,
SEC_TYPE_DESC,
MF_IND,MONEY_MARKET_IND,SHARE_CLASS,CURR_PRICE,CURR_PRICE_DT,PAY_FREQ_CD,COUPON_RT,MATURITY_DT
,MODDY_RT,SP_RT,CALL_PRICE,DIVIDEND,FACTOR_DT,FACTOR,NOTE,INTEGRATION_TSTP
FROM STAGE.W_BDA_SECURITY_DM_4

INSERT INTO SAMPLE_SCHEMA.W_BDA_SECURITY_DM_Dist
(ROW_ID,
SEC_NO,
CUSIP,
SEC_LONG_NAME,
SEC_SHORT_NAME,
SEC_SYMBOL,
SEC_CLASS_CD,
OLD_SEC_TYPE_CD,
NEW_SEC_TYPE_CD,
SEC_TYPE_DESC,
MF_IND,MONEY_MARKET_IND,SHARE_CLASS,CURR_PRICE,CURR_PRICE_DT,PAY_FREQ_CD,COUPON_RT,MATURITY_DT,
MODDY_RT,SP_RT,CALL_PRICE,DIVIDEND,FACTOR_DT,FACTOR,NOTE,INTEGRATION_TSTP
)
SELECT
B.NEW_ROW_NUM, B.SEC_NO, B.CUSIP, B.SEC_LONG_NAME,
 B.SEC_SHORT_NAME, B.SEC_SYMBOL,
 B.SEC_CLASS_CD, B.OLD_SEC_TYPE_CD, B.NEW_SEC_TYPE_CD, B.SEC_TYPE_DESC, B.MF_IND, B.MONEY_MARKET_IND,
 B.SHARE_CLASS, B.CURR_PRICE, B.CURR_PRICE_DT, B.PAY_FREQ_CD, B.COUPON_RT, B.MATURITY_DT,
 B.MODDY_RT, B.SP_RT, B.CALL_PRICE, B.DIVIDEND, B.FACTOR_DT, B.FACTOR, B.NOTE, B.INTEGRATION_TSTP
 FROM STAGE.W_BDA_SECURITY_DM_4_ROWNUM B
 END
```
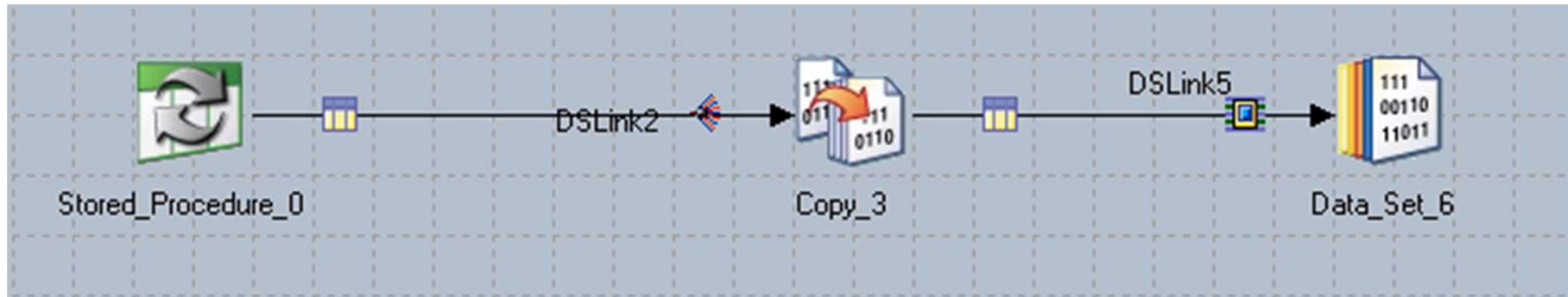
Stage Table with sequencer

Target Table

Data Migration Jumpstart Engineering

# DataStage – Job to Execute Stored Procedure

# Stored Proc General Setup

# Set Up Syntax for Stored Proc Exec

---

Azure Synapse Stored Procedure Name

T-SQL Stored Procedure Call

**Stored_Procedure_0 - STPPX stage**

Stage | Output

Stage name:
Stored_Procedure_0

Data Connection:

General | Syntax | Parameters | Error Codes | NLS Map | Advanced

Procedure name:
dbo.STAGE_INSERT

Procedure type:
Source

Database procedure type:
Stored procedure

Procedure call syntax:

☐ Generate procedure call

EXEC dbo.STAGE_INSERT;

OK    Cancel    Help

◄ ► ►| \ Stored_Proc_Check_1 /

# Set up "Copy" process Input Properties

# Set up "Copy" Process Output Properties for Stored Proc

Set up "Copy" process for Output Mapping

Set up "Copy" process Output Columns

# Large Data Ingestion in Azure Synapse Strategy

- Large Data Set ingestion should follow a pattern
  - Land data intro Azure Blob Container or Data Lake Store
  - Ingest into a staging table in Azure Synapse Analytics using COPY INTO
  - CTAS or Insert into Final Table
  - Refresh statistics as needed
- DataStage can orchestrate the reading from source (e.g. flat files) and ingesting into Blob Storage as a first step

Azure Blob / Data Lake Store  —— COPY INTO ——>  Azure Synapse Analytics SQL Pool

Data Migration Jumpstart Engineering

# Writing Data Into Azure Storage with DataStage

- Large Data Set ingestion can take advantage of COPY INTO process from Azure Synapse Analytics from data landed in Blob Storage
- Requirements: IBM Infosphere Information Server Datastage 11.7fp1 and above
- High Level Steps:
  1. Configure Azure Storage Connector Connection Properties
  2. Configure Azure Storage Connector to write to Azure Blob Storage
  3. Additional Configuration for Parallel Write
- [Documentation Link](#)



input_db2_data — data — azure_storage_file

Data Migration Jumpstart Engineering

# Bulk Loading into Azure Synapse with DataStage

- COPY INTO statement can be executed:
  - Within a Stored Procedure
  - Or as an ad-hoc T-SQL
- COPY INTO will take the data from Blob Storage/ADLS and ingest into Azure Synapse Analytics in parallel
- Data should be landed on an uncompressed rowstore table for performance reasons
- T-SQL should execute a last step of ingesting or CTAS into a final table.
- Recompute of statistics is always recommended when data has changed significantly

# Bulk Loading with COPY INTO statement

```
CREATE SCHEMA stage;
-- Create the Staging Table
-- Drop Table [stage].[SalesData]
Create table [stage].[SalesData](
                Region varchar(50)
                ,Product varchar(50)
                ,SaleDate date
                ,Amount DECIMAL (20,10)
                )
with (distribution = round_robin, HEAP);


/* Documentation
https://docs.microsoft.com/en-us/sql/t-sql/statements/copy-into-transact-sql?view=azure-sqldw-latest
*/

COPY INTO [stage].[SalesData] (Region 1, Product 2, SaleDate 3, Amount 4)
FROM 'https://acccount.blob.core.windows.net/yourcontainer'
WITH (
    FILE_TYPE = 'CSV',
    CREDENTIAL=(IDENTITY= 'Shared Access Signature', SECRET='Your secret here'),
                FIELDQUOTE = '"',
                FIELDTERMINATOR=',',
                ROWTERMINATOR = '\n',
                Identity_INSERT = 'OFF',
                ENCODING = 'UTF8',
                DATEFORMAT = 'ymd',
                FIRSTROW = 2
);

CREATE SCHEMA Sales;
-- CTAS with HASH distribution on Date and Columnar format
-- DROP TABLE Sales.Transactions
CREATE TABLE Sales.Transactions
WITH
(
 DISTRIBUTION = HASH(SaleDate)
 ,CLUSTERED COLUMNSTORE INDEX
) AS SELECT * FROM [stage].[SalesData];

Create Statistics stat_stage_salesdata_Region on Sales.Transactions (Region) with fullscan;
Create Statistics stat_stage_salesdata_Product on Sales.Transactions (Product) with fullscan;
Create Statistics stat_stage_salesdata_SaleDate on Sales.Transactions (SaleDate) with fullscan;
Create Statistics stat_stage_salesdata_Amount on Sales.Transactions (Amount) with fullscan;
```

No compression

Not necessary if using AAD / Role based access control

CTAS into final table, with compression and distribution

Statistics are created

Data Migration Jumpstart Engineering