

DEVELOPING DATA PIPELINES WITH AZURE SYNAPSE

By Temi | November 29, 2023

Agenda

Content

Page

- | | |
|--|---|
| 1. Big Data Integration with Azure Synapse | 3 |
| 2. Synapse Features | 4 |
| 3. Pipeline Components | 5 |
| 4. ETL Pipeline Demo – Developer Productivity Tracking | 8 |

Big Data Integration with Azure Synapse



Azure Synapse Features

1 Linked Services

- Stores connection information to external datasources or compute like storage account, Azure SQL DB (postgres, mysql), Azure Keyvault

2 Integration Datasets

- Stores connection information to datasets themselves

3 Integration Runtime

- It is referenced by a linked service or a pipeline activity
- Provides compute where pipeline activity is dispatched
- **Types:** Azure or Selfhosted

4 Version Control with Git

- Git integration with Github or Azure pipelines

5 Database Objects

- Schemas
- Tables and Views
- Store procedures

6 Synapse Pipelines and activities

- Orchestrate ETL pipelines using components / activities

Pipeline Components – Copy Data Activity

Copy data from one location to another

- From SQL Server database to Storage account in Parquet format.
- Copy text files (CSV) format from an on-premises file system and write Storage account in Avro format.
- Copy data from Storage account to SQL Pool

Copy data



LoadStagingTableFromParquet

Pipeline Components – Notebook Activity

Run a Notebook on a Spark Pool

- Process data in Storage account using Pyspark and Pandas
- Save data to storage account using the abfss protocol (managed identity of synapse requires Storage Blob Contributor IAM permission)

Notebook



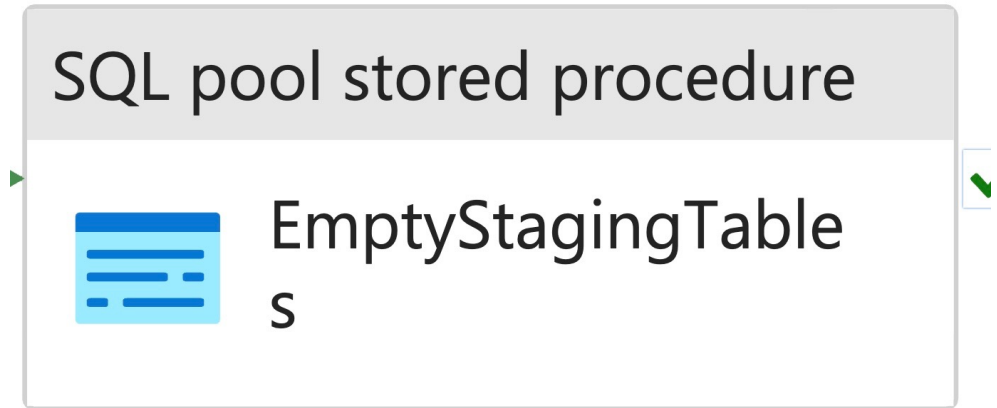
ProcessJiraGitData



Pipeline Components – Stored Procedure Activity

Run a stored procedure in the Dedicated SQL pool

- Execute stored procedures written in TSQL in an ETL pipeline



ETL Pipeline Demo – Developer Productivity Tracking

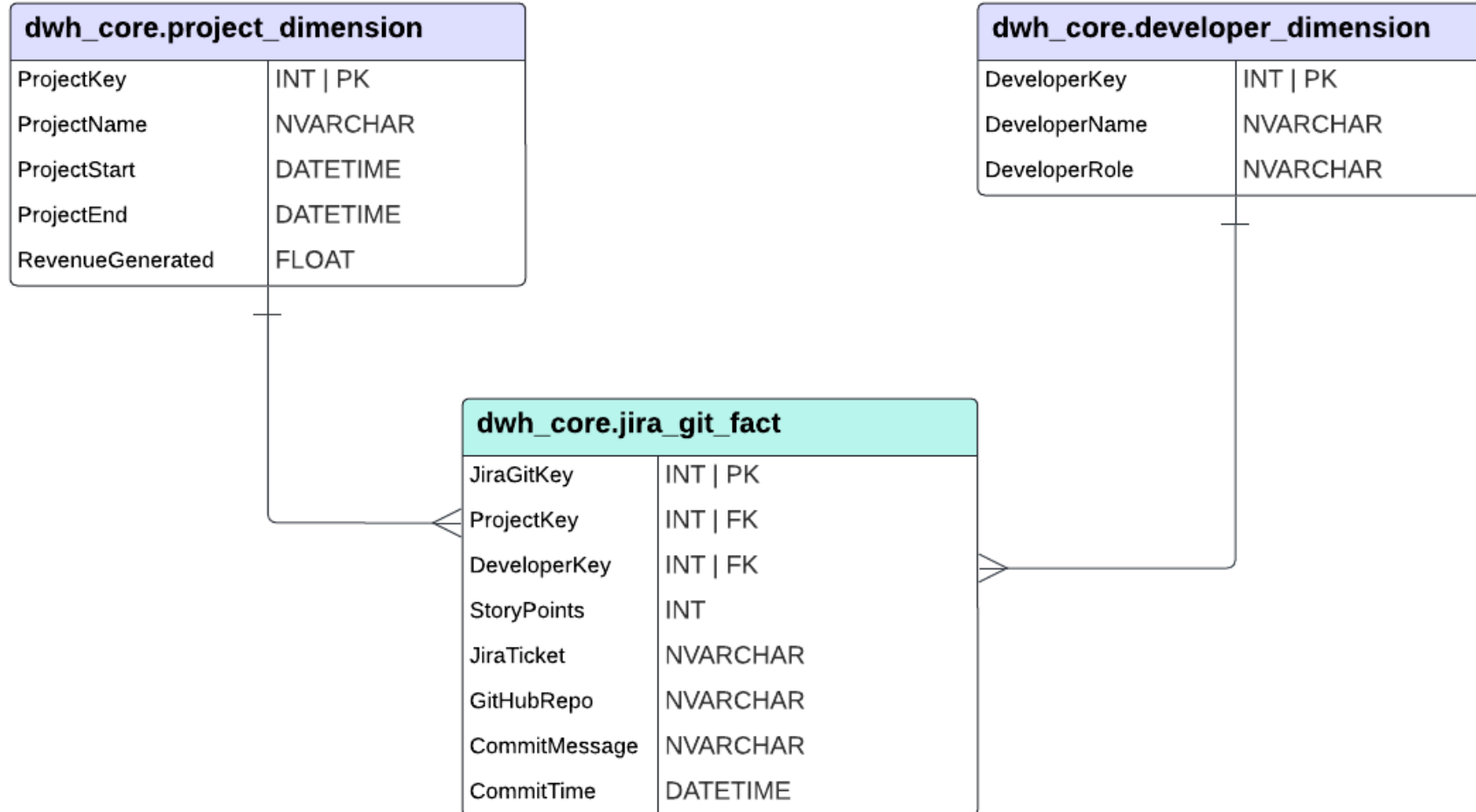
Scenario

External data from Git and Jira is integrated into our Datawarehouse in Synapse.

Given information on our projects, developers and their activity on Github and Jira, KPIs are created to measure their productivity.

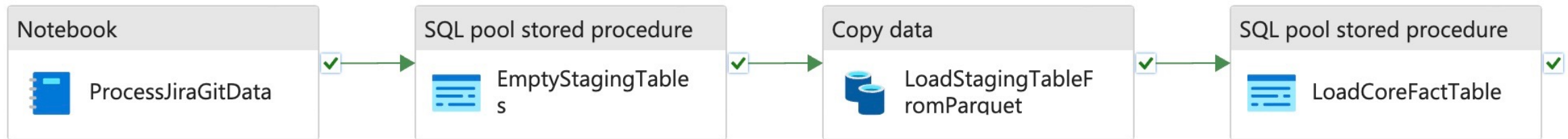


ERD Diagram



Developer Productivity Pipeline for Jira Git Table

1. **Notebook:** process the `jira_git_fact.csv` into `jira_git_fact.parquet`
2. **CopyData:** copy from `jira_git_fact.parquet` into staging table – `dwh_staging.jira_git_fact`
3. **Stored Procedure:** copy from `dwh_staging.jira_git_fact` into core table - `dwh_core.jira_git_fact`



Load CSV Pipeline for Developer and Project Tables

1. **CopyData**: copy from **developer_dimension.csv** into staging table – **dwh_staging.developer_dimension**
2. **Stored Procedure**: copy from **dwh_staging.developer_dimension** into core table – **dwh_core.developer_dimension**



Load CSV Pipeline Parameters

PARAMETER NAME	DEVELOPER_DIM_VALUES	PROJECT DIM VALUES
CsvFolder	bronze/developer_dim	bronze/project_dim
CsvFileName	developer_dimension.csv	project_dimension.csv
StagingStoredProcedure	dwh.empty_staging_tables	dwh.empty_staging_tables
CoreStoredProcedure	dwh.load_core_developer	dwh.load_core_project
StagingTableName	developer_dimension	project_dimension

THANK YOU