# A DATA ENGINEERS TOOLKIT

# SO WHAT IS DATA ENGINEERING?

# THE OLD WORLD



Enterprise Data Warehouse

Source Systems

Stage     Clean     Warehouse

# PACKAGE ORCHESTRATION
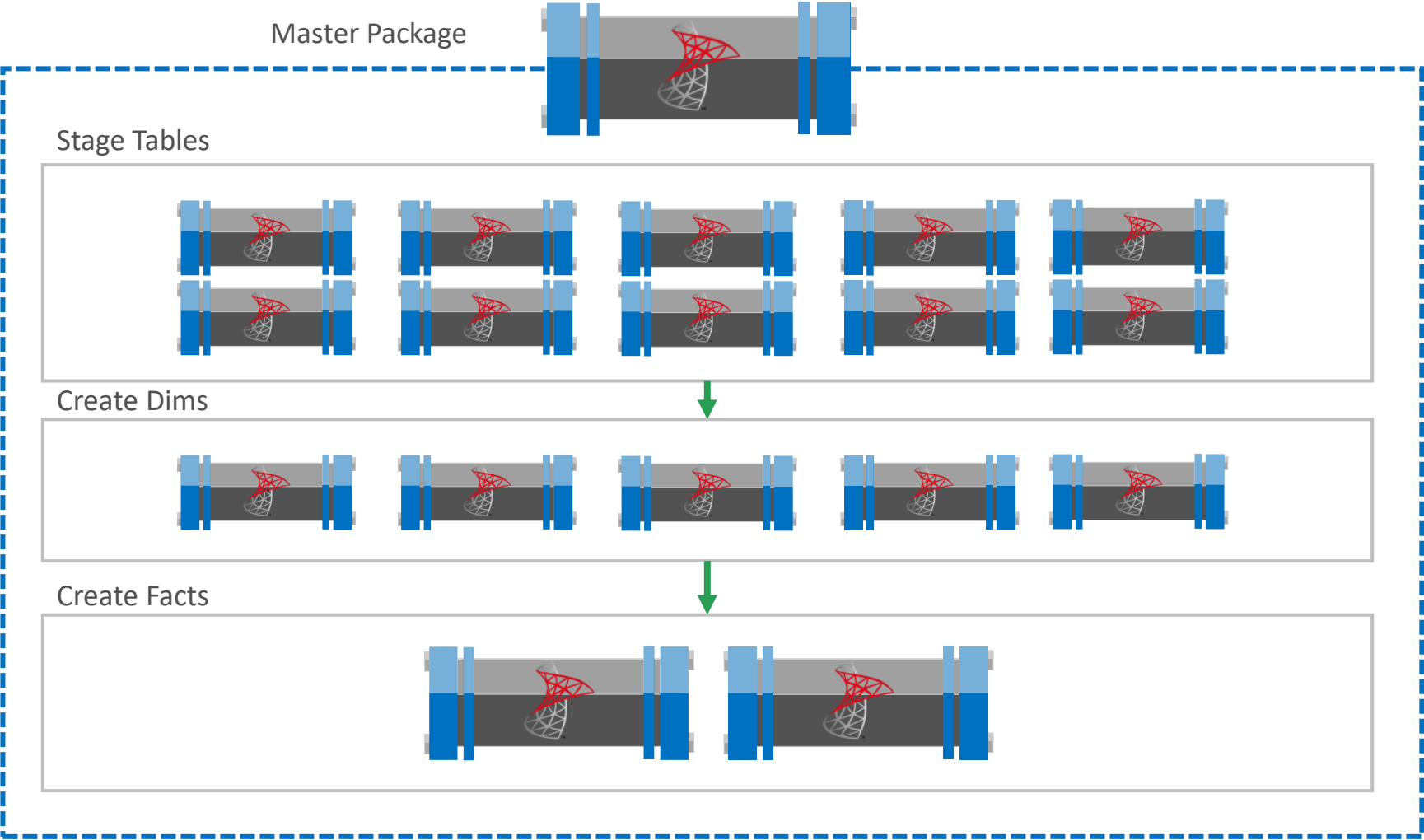
# OUR ETL TOOL CHECKLIST

- Low Dev Effort
- Many Use Cases
- Flexible Formats
- Elastic Scaling
- Agile
- Support

FAILED

# SO MOVE TO AZURE, IT'S EASY RIGHT?

## DATABASES (19)

- Azure Cosmos DB
- Azure SQL
- SQL databases
- Azure Database for MySQL servers
- Azure Database for PostgreSQL servers
- Azure Database for MariaDB servers
- SQL servers
- Dedicated SQL pools (formerly SQL DW)
- Azure Database Migration Services
- Azure Cache for Redis
- SQL Server stretch databases
- Data factories
- SQL elastic pools
- Virtual clusters
- Managed databases
- Elastic Job agents     PREVIEW
- SQL managed instances
- SQL virtual machines
- SQL Server registries     PREVIEW

## ANALYTICS (14)

- Dedicated SQL pools (formerly SQL DW)
- Azure Databricks
- HDInsight clusters
- Data factories
- Power BI Embedded
- Stream Analytics jobs
- Data Lake Analytics
- Analysis Services
- Event Hubs
- Event Hubs Clusters
- Log Analytics workspaces
- Data Lake Storage Gen1
- Azure Data Explorer Clusters
- Power Platform     PREVIEW

# MORE REALISTICALLY, WE HAVE SEVERAL KEY OPTIONS

**Azure Data Factory**

Visual cloud-native Orchestration. Excellent at moving data between sources

**Databricks**

Spark as a service, powerful multi-language big data engine. Supports SQL but best with Python/Scala

**Azure Synapse Analytics**

Suite of tools ranging from orchestration to engineering & full data warehousing

# A MODERN ARCHITECTURE

Orchestration

Azure Databricks

RAW | BASE

DELTA

Data Factory

Data Sources

Data Lake Store Gen 2

Azure Synapse Analytics

PowerBI

# ORCHESTRATION & INTEGRATION

# AN ADF PRIMER



Copy Data

Transform

Performs copy activities using it's own scalable compute

**Azure Data Factory**

Triggers other components for heavy lifting & intensive transformation

JSON

ADVANCING ANALYTICS

Sources

Copy Data

Transform Activities

Destinations

ADVANCING ANALYTICS

# DEVELOPER TOOLS

# MONITORING

📅 **Custom Range** 04/02/2019 8:00 AM - 04/07/2019 8:00 AM ⌄    🌐 **Time Zone** (UTC+00:00) Dublin, Edinburgh, Li... ⌄    ⚪ View All Rerun History    ▽② Filter

**All**   Succeeded   In Progress   Queued   Failed   Cancelled

| ☐ | Pipeline Name ▽ | Actions | Run Start ↕ | Duration | Triggered By | Status | Parameters | Annotations ▽ | Error | RunID |
|---|---|---|---|---|---|---|---|---|---|---|
| | RunNotebooks | ▶ ↻ | 04/03/2019, 6:08:10 PM | 00:05:49 | Manual trigger | ✅ Succeeded | | | | 13ca5395-bb39-4559-a14 |
| | RunNotebooks | ▶ ↻ | 04/03/2019, 4:03:46 PM | 00:00:37 | Manual trigger | ✅ Succeeded | | | | 5b101cf8-e105-4e1c-8f82 |

📅 **Custom Range** 04/02/2019 8:00 AM - 04/07/2019 8:00 AM ⌄    🌐 **Time Zone** (UTC+00:00) Dublin, Edinburgh, Li... ⌄

**⫘ Pipeline**                                                                    **💼 Activity**

**100%**
SUCCEEDED RUNS

SUCCEEDED RUNS
**2**

SUCCEEDED RUNS
**2**

**ADVANCING
ANALYTICS**

# A QUICK LOOK AT ADF

- The ADF Studio
- Object Concepts

# SSIS – LIFT & SHIFT POTENTIAL

Pipeline

SSIS Activity

SSIS Catalog

SSIS IR

ADVANCING ANALYTICS

# DATA FACTORY RECAP

- Orchestrates all data workflows in Azure
- Best method of onboarding data to Azure
- Use parameters, forEach and child executions

- Low Dev Effort
- Many Use Cases
- Flexible Formats
- Elastic Scaling
- Agile
- Supportable

PASSED

ADVANCING ANALYTICS

COMPUTE APPROACHES IN AZURE

# QUICK SPARK OVERVIEW

Spark is a distributed, scalable data processing engine.

It can query **structured** and **non-structured** data

You can use **Python**, **Scala**, **R**, **C#** or **SQL** to interact with it

CSV

PYTHON

ADVANCING
ANALYTICS

## DataFrame

- Schema
- Format
- Location

```
df = (spark
        .read
        .schema(newSchema)
        .format(fileFormat)
        .load(dataLocation)
     )
```

ADVANCING
ANALYTICS

# DISTRIBUTED COMPUTE

# DISTRIBUTED COMPUTE

# DISTRIBUTED COMPUTE

# DISTRIBUTED COMPUTE

# AZURE DATA FACTORY – MAPPING DATA FLOWS

# AZURE DATA FACTORY – MAPPING DATA FLOWS

Mapping Data Flows are a **GUI-Based transformation tool**, modelled on SSIS Data Flows

They use a **managed Spark Engine** which gives you the power and scale of Spark, but without the programming overhead

They are **not as flexible** as full spark but are getting better all the time!

ADVANCING
ANALYTICS

# AZURE DATABRICKS

Microsoft Azure | Databricks        Portal    simon@advancinganalytics.co.uk

## Dynamic Validation (Python)

AdventureWorks Table : Product

### Read the schema json for our selected file

I've stored a schema file for each of the data files in my lake. I can pick up the right file for the dataset selected by my widget

Cmd 7

```python
1   #Load the relevant libraries to build schemas and read JSON
2   from pyspark.sql.types import *
3   import json
4
5   #Inject our filename into the lake path
6   schemaLocation = f"/mnt/dblake/RAW/Public/Adventureworks/SalesLT.{fileName}.json"
7
8   #Read the json file contents
9   jschemadf = sqlContext.read.text(schemaLocation)
10
11  #Pull out the first value (it's all one value but the reader turns it into a dataframe)
12  jschema = jschemadf.first().value
13
14  #Convert our JSON schema into a pyspark Struct which can be applied directly to a dataframe
15  newSchema = StructType.fromJson(json.loads(jschema))
16  newSchema
```

▶ (1) Spark Jobs

▶ 📄 jschemadf: pyspark.sql.dataframe.DataFrame = [value: string]

Out[17]: StructType(List(StructField(ProductID,IntegerType,true),StructField(Name,StringType,true),StructField(ProductNumber,StringType,true),StructField(Color,StringType,true),StructField(StandardCost,StringType,true),StructField(ListPrice,DoubleType,true),StructField(Size,StringType,true),StructField(Weight,StringType,true),StructField(ProductCategoryID,StringType,true),StructField(ProductModelID,IntegerType,true),StructField(SellStartDate,StringType,true),StructField(SellEndDate,StringType,true),StructField(DiscontinuedDate,StringType,true),StructField(ThumbNailPhoto,StringType,true),StructField(ThumbnailPhotoFileName,StringType,true),StructField(rowguid,StringType,true),StructField(ModifiedDate,StringType,true)))

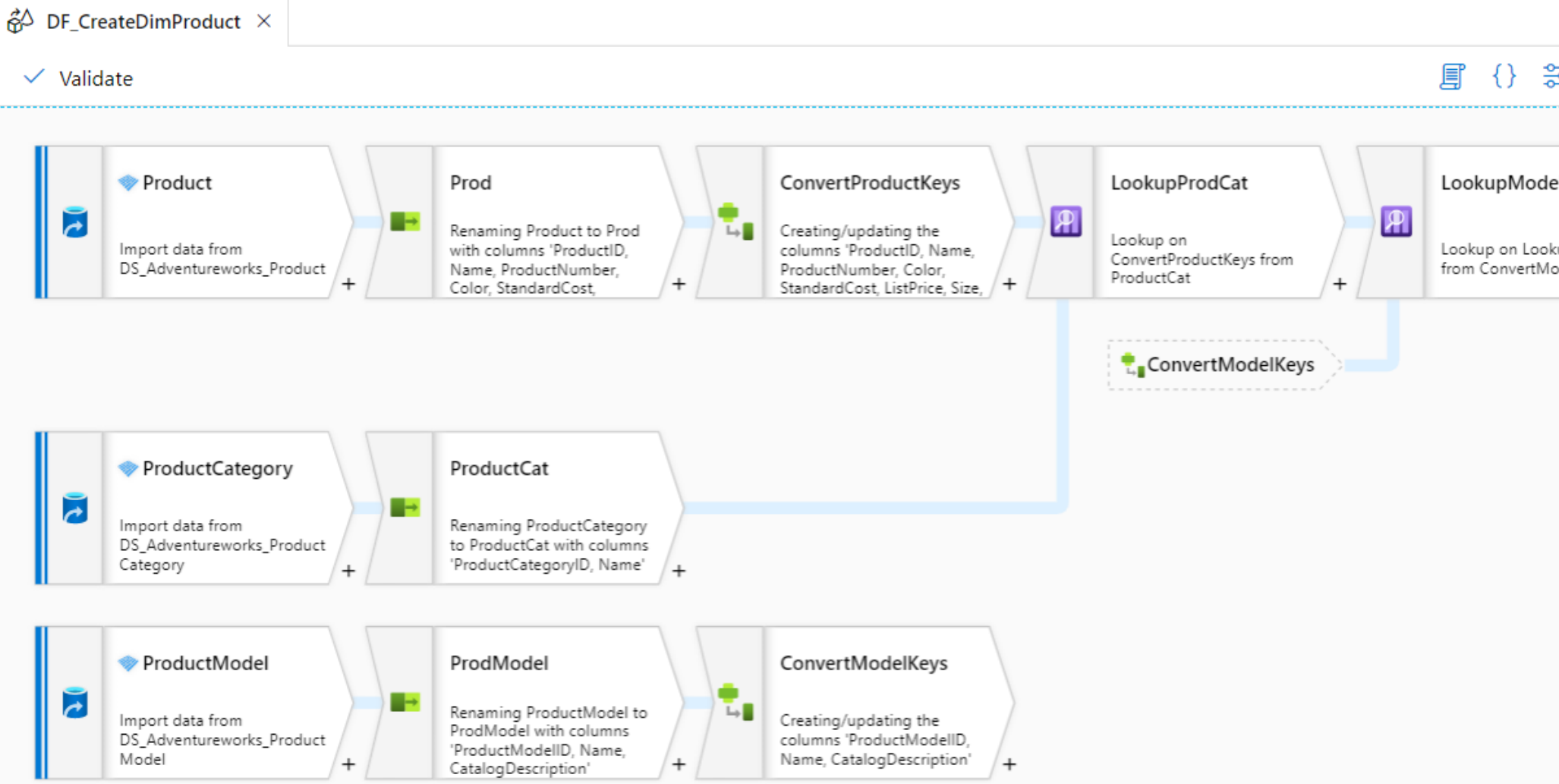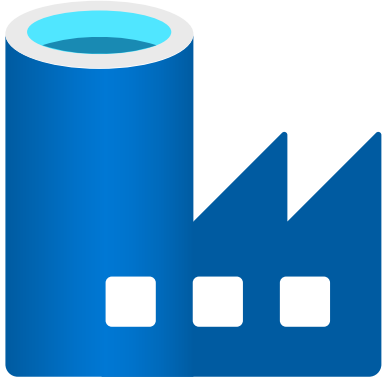Command took 0.38 seconds -- by simon@advancinganalytics.co.uk at 26/09/2020, 15:33:57 on Runtime7

Cmd 8

### We have a schema, now we need to create a dataframe

We can derive the path of our dataset in the same way as we did with the schema. We then combine schema and data location in a new dataframe

We're also going to use "_corrupt_record", this is a system field which will only be populated if a row fails to parse into the structure we've provided

Cmd 9

ADVANCING
ANALYTICS

# AZURE DATABRICKS

Databricks is a third-party company, founded by the **team who invented spark**

They provide an **Azure-native, managed Spark platform**

Databricks will generally have the **most advanced spark engine** and maintain a fast release cadence

DATABRICKS SQL ANALYTICS

Analytics

DS/ML

SQL Analytics

Data Science Workspace & ML Workflow

Delta Lake

Operational and External Structured, Semi-Structured, and Unstructured Data

ADVANCING ANALYTICS

Where to start...

# AZURE SYNAPSE ANALYTICS

**Data Pipelines**

**Mapping Data Flows**

**Azure Synapse Studio**

**Dedicated SQL Pools** *(SQLDW)*

**Serverless SQL Pools**

**Provisioned Spark Pools**

**Monitoring**

**Data Lake Store Gen 2**

**Metadata Store**

**Management**

ADVANCING
ANALYTICS

- Billed Per Session Uptime
- Scala, Python, C#, SQL
- Dynamic Workflows, Machine Learning & Unusual Data Types

- Session management is... interesting

ADVANCING ANALYTICS

# SERVERLESS SQL POOLS

- Billed Per TB Read
- T-SQL
- Ad-hoc/Occasional access

- Unpredictable Billing
- Very Black-box (but is that bad?)

ADVANCING ANALYTICS

# DEDICATED SQL POOLS

- Billed Per Hour
- T-SQL
- Huge Datasets & Formal Modelling

- Inflexible Scaling
- Can be complex to distribute tables

ADVANCING
ANALYTICS

# AZURE SYNAPSE ANALYTICS

- Synapse Workspace Overview

- Serverless SQL

- Spark Pools

A MODERN ARCHITECTURE

Orchestration

Azure Databricks

RAW | BASE

DELTA

Data Factory

Data Sources

Data Lake Store Gen 2

Azure Synapse Analytics

PowerBI

# MODERN DATA WAREHOUSES

Azure SQL
Datawarehouse

RAW   BASE   ENRICHED   CURATED

Manual
Upload

Data Sources

SQL