



Stuart Tomlinson
Microsoft Cloud Consultant

Discover Big Data with Azure Data Factory



Agenda

What is Azure Data Factory

Data Integration

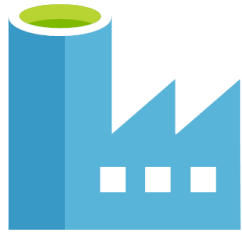
Data Factory Components

Demonstration



Big Data in the Cloud: **What is Data Factory?**





What is **Azure Data Factory**?

2016- Introduced into Azure Preview (V1)
And was initially marketed as a serverless SSIS platform in the Azure cloud.

2020- Microsoft introduced ADF V2 with many more capabilities, with code free dataflows with over 90 connectors for data.

Azure Data Factory now Natively Integrates SQL, HD Insights, Oracle, Data Bricks and Synapse Big Data tools and many more, while still allowing SSIS packages to be migrated into ADF.

WHAT IS AZURE DATA FACTORY

READ THE DOCS

A CLOUD-BASED DATA INTEGRATION SERVICE ORCHESTRATES DATA MOVEMENT & TRANSFORMATION BETWEEN DIVERSE DATA SOURCES & CLOUD COMPUTE RESOURCES AT SCALE

DATA INTEGRATION CHALLENGE

WHAT PROBLEM DOES IT SOLVE?

BIG DATA: DIVERSE DATA SOURCES, DURATIONS, STORED DATA, CLEANING & ENRICHED, ACTIONABLE INSIGHTS, DASHBOARDS, REPORTS

DATA STREAMS, DEVICES, DATA INGEST, DATA TRANSFORMATION, DATA ANALYSTS

HELP! We need to collect all this data

DATA OPS TEAM

DATA MOVEMENT

DATA ANALYSTS

WE NEED DATA WITH RELEVANT CONTEXT TO DO OUR ANALYSIS

Raw, unorganized data stored in relational, non-relational, and other storage systems

HOW CAN WE ORCHESTRATE (DATA MOVEMENT) AND OPERATIONALIZE (WORKFLOW)

AZURE DATA FACTORY

We need actionable business insights but raw data lacks the correlating content (by itself) for analysis

Big data needs a scalable service than can orchestrate the data movement and operationalize the data processing workflows!

THE SCENARIO YOU RECOGNIZE

GAME DEV COMPANY

HAVE: Petabytes of game log data in the cloud

WANT: Actionable insight into my customers

US gaming preferences, US player demographics, US user behaviors

USES

- Develop new features
- Improve player experience
- Upsell, cross-sell and drive business growth

WHAT I NEED

- reference data from our on-prem services
- customer info
- marketing campaign game metadata

* combined with gameplay data logs from cloud

game play behaviors, success, failure rates etc.

WHAT I MUST DO

- collect data from sources
- process data to transform it
- publish data into warehouse for analysts

AUTOMATE THIS WORKFLOW

TRIGGER IT BY AN EVENT OR MANUALLY

MONITOR & MANAGE DAILY SCHEDULES

HOW DOES ADF WORK?

LET'S BREAK IT DOWN

CONNECT & COLLECT

- Build an information production system to connect all data sources and adapt ingest to their diverse intervals and speeds!
- Collect all data in a centralized location to facilitate processing (e.g. transformation) next

TRANSFORM & ENRICH

- Transform the collected data using actions to aggregate, filter, clean etc. Use code-free UI based mapping data transformation graph creation or use compute to transform data by hand

CI/CD PUBLISH

- Manage data pipeline ops using Azure DevOps and Github! Incrementally develop and deploy your ETL process then publish processed data into Azure Data Warehouse, Azure SQL database, Azure Cosmos DB - or your BI analysis engine

MONITOR & ALERT

Use Azure Monitor, APZ, PowerShell, health panels on Azure Portal!

LET'S TALK DATA INTEGRATION PATTERNS

The most common pattern is ETL

EXTRACT = connect to sources, copy data to central store

TRANSFORM = process the data for analysis, data transformation at scale

LOAD = move the data to data warehouse or analytics engines for business insights

Another pattern is ELT where raw (native) data is itself loaded (stored) before the transformation phase...

AZURE DATA FACTORY HAS CODE-FREE ETL AS A SERVICE!

You always have the option to do hand-coded transformation using Azure compute... but mapping data flows provide a UI-based wizard to simplify your pipeline setup..

5 ASPECTS IT COVERS!

- INGEST DATA
- CONTROL FLOW
- DATA FLOW
- SCHEDULE OPS
- MONITOR OPS

7 THINGS TO KNOW ABOUT AZURE DATA FACTORY

- ENTERPRISE READY** Data Integration At cloud scale!
- ENTERPRISE DATA READY** 90+ connectors! It just works.
- CODE-FREE TRANSFORMATION** UI driven mapping data flows
- RUN CODE ON ANY AZURE COMPUTE** For hands-on data transformations
- MANY SSIS PACKAGES RUN ON AZURE** (Migrate on-prem SSIS in 3 steps)
- ADF CAN MAKE DATA OPS SEAMLESS** Source control, automated deploys, simple templates
- SECURE DATA INTEGRATION** Managed virtual networks protect against data exfiltration, simplify "your" networking!

ADDITIONAL CONCEPTS

1 CONTROL FLOW

REFERS TO THE ORCHESTRATION OF PIPELINE ACTIVITIES INCLUDING:

- CHAINING ACTIVITIES
- BRANCHING ACTIVITIES
- PASSING ARGUMENTS (for pipeline run)
- CUSTOM STATE PASSING
- LOOPING CONTAINERS

2 PIPELINE RUN

AN INSTANCE OF PIPELINE EXECUTION

INSTANTIATE PIPELINE BY PASSING ARGUMENTS (VALUES) TO PARAMETERS (PLACEHOLDERS) DEFINED IN IT

PASS ARGUMENTS MANUALLY OR BY DEFINING TRIGGERS

3 TRIGGERS

UNIT OF PROCESSING THAT DETERMINES WHEN TO KICK OFF PIPELINE RUNS.

DIFFERENT TRIGGER TYPES FOR VARIOUS EVENTS!

4 PARAMETERS

"READ-ONLY" KEY/VALUE PAIRS POPULATED FROM RUN CONTEXT AT EXECUTION

BOTH "DATASET" AND "LINKED SERVICE" ARE STRONGLY-TYPED, REUSABLE, REFERENCEABLE PARAMETER ENTITIES!!

ACTIVITY REFERENCES DATASET TO CONSUME DATA PROPERTIES, AND LINKED SERVICE FOR THE CONNECTION INFO (to data store for copy, compute resource for flow)

5 VARIABLES

ARE USED INSIDE PIPELINES TO STORE TEMPORARY VALUES, STATE

USED WITH PARAMETERS TO PASS VALUES BETWEEN ACTIVITIES, DATA FLOWS, PIPELINES

6 INTEGRATION RUNTIME (IR)

INTEGRATED RUNTIME (IR) IS COMPUTE INFRASTRUCTURE USED BY ADF TO PROVIDE FULLY MANAGED

- DATA FLOWS (transformation)
- DATA MOVEMENT (movement)
- ACTIVITY DISPATCH (route to service)
- SSIS PACKAGE EXECUTION (in managed)

CREATE & MANAGE DATA TRANSFORMATION GRAPHS THAT WORK ON ANY SIZE DATA

BUILD UP REUSABLE LIBRARY OF DATA TRANSFORMATION ROUTINES!

EXECUTE TRANSFORM ROUTINES IN PIPELINES (AS ACTIVITY) TO SCALE!

HEY PRESTO!

ADF EXECUTES LOGIC ON SPARK CLUSTER FOR YOUR DATA (NO NEED TO MANAGE OR MAINTAIN YOUR OWN CLUSTER)

ACTIVITIES

REPRESENT A SINGLE PROCESSING STEP IN THE PIPELINE

EXAMPLES:

- COPY ACTIVITY: copy data from one store (source) to other (destination)
- HIVE ACTIVITY: run a Hive query on an HDInsight cluster

THREE TYPES OF ACTIVITY SUPPORTED BY ADF

- Data Movement
- Data Transformation
- Control

MAPPING DATA FLOWS

4 LINKED SERVICES

TWO USES

- REPRESENT DATA SOURCE
- REPRESENT COMPUTE RESOURCE

REPRESENT CONNECTION INFORMATION TO LINK ADF TO EXTERNAL SERVICES

WORKS CLOSELY WITH DATASETS. THIS SPECIFIES HOW TO CONNECT & DATASET DEFINES WHAT DATA LOOKS LIKE...

DATA STORE

DATASET (VIEW)

PIPELINES

A PIPELINE IS A LOGICAL GROUPING OF ACTIVITIES THAT PERFORM ONE UNIT OF WORK

AN ADF INSTANCE CAN HAVE ONE OR MORE PIPELINES!

Example: Pipeline that ingests data from Azure Blob (A1), runs Hive query on HDInsight to partition it (A2) and moves result to next store (A3)

Benefit: Manage correlated activities as a single entity!

Execution: Activities may be chained (run sequentially) or be independent (run in parallel) within pipeline!

HERE ARE THE KEY TERMS AND CONCEPTS YOU NEED TO KNOW TO USE ADF!

AND A FEW KEY TERMS

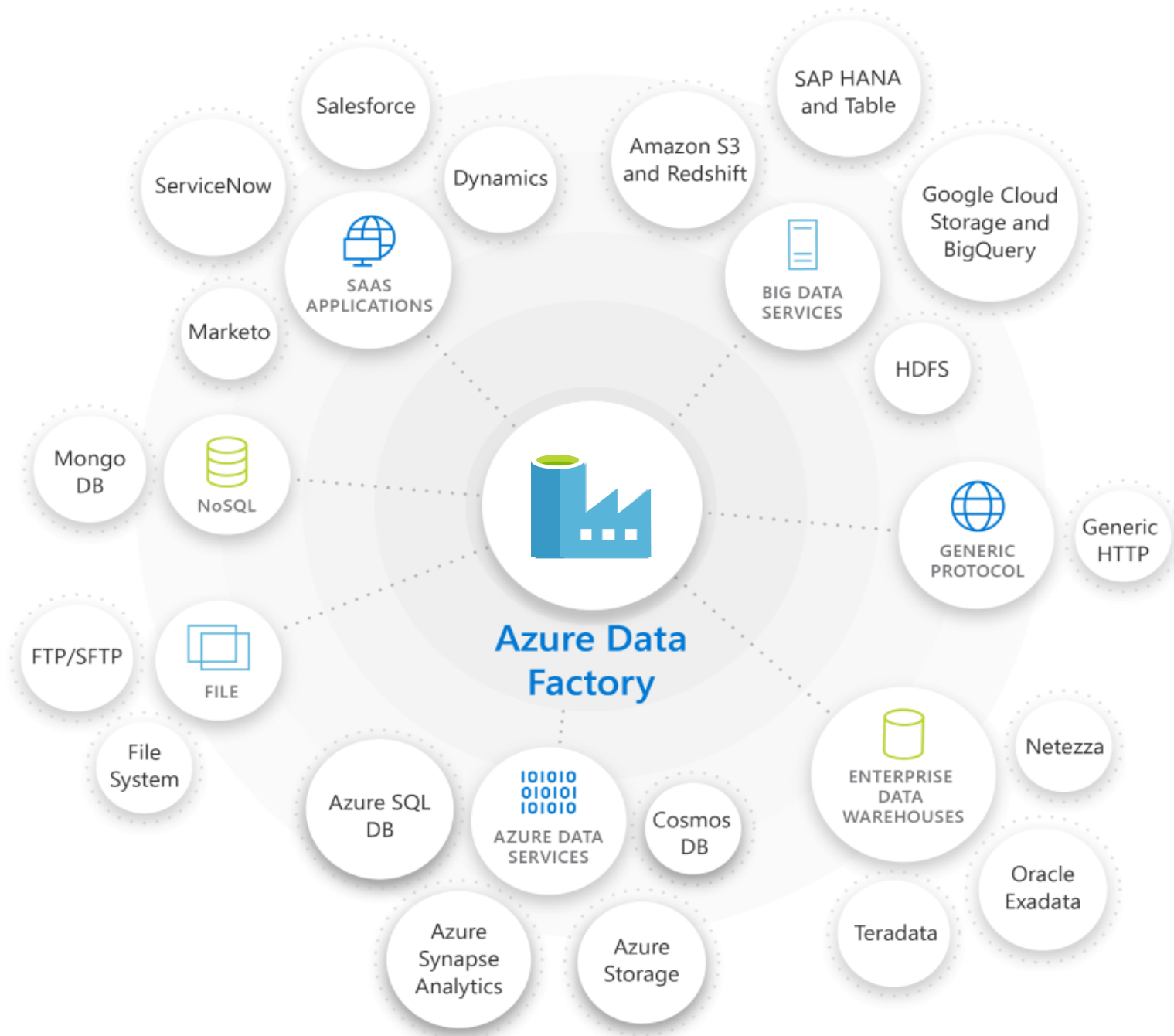
- CONTROL FLOW
- PIPELINE RUN
- TRIGGERS
- PARAMETERS
- VARIABLES

ADF TOOLKIT

LET'S TALK ABOUT COMPONENTS OF AZURE DATA FACTORY

6 MAIN CONCEPTS

- PIPELINES
- ACTIVITIES
- DATASETS
- LINKED SERVICES
- DATA FLOWS
- INTEGRATION RUNTIMES



INGEST



- Multi-cloud and on-prem hybrid copy data
- 90+ native connectors
- Serverless and auto-scale
- Use wizard for quick copy jobs

With over 90 native connectors ADF can ingest and transform Data from On Premise File shares, AWS containers, Azure Storage, Oracle databases, SQL and many more.

Data can then be copied and transformed into Azure and on Premise SQL databases , Files shares and even AWS databases for additional reporting or transformation while leaving the source data intact.





Code-Free ETL as a Service

INGEST



- Multi-cloud and on-prem hybrid copy data
- 90+ native connectors
- Serverless and auto-scale
- Use wizard for quick copy jobs

CONTROL FLOW



- Design code-free data pipelines
- Generate pipelines via SDK
- Utilize workflow constructs: loops, branches, conditional execution, variables, parameters, ...

DATA FLOW



- Code-free data transformations that execute in Spark
- Scale-out with Azure Integration Runtimes
- Generate data flows via SDK
- Designers for data engineers and data analysts

SCHEDULE



- Build and maintain operational schedules for your data pipelines
- Wall clock, event-based, tumbling windows, chained

MONITOR



- View active executions and pipeline history
- Detail activity and data flow executions
- Establish alerts and notifications

Data Factory allows us to create data-driven workflows in the cloud for orchestrating and automating data movement and data transformation.

Semi-Structured Data



Excel Files



Power Apps



SharePoint

Ingest



Azure Data
Factory

Store (New)



Azure SQL
Relational
Database

Model (New)



Power BI
Datasets

Serve

Dashboards, Reports, Alerts



Monitoring / Analysis / Review

Store (Existing)



DataLake



Other
Databases



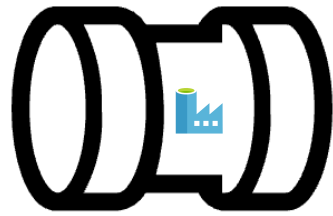


Data Factory **Components?**

Azure Data Factory is composed of these key components

- Pipelines
- Activities
- Datasets
- Linked services
- Data Flows
- Integration Runtimes

These components work together to provide the platform on which you can compose data-driven workflows to move and transform data.



ADF Pipeline

A pipeline is a logical grouping of activities that performs a unit of work. Together, the activities in a pipeline perform a task.

The benefit of this is that the pipeline allows you to manage the activities in a set instead of managing each one individually.

The activities in a pipeline can be chained together to operate sequentially, or they can operate independently in parallel all in the same pipeline run.

.

CONTROL FLOW



- Design code-free data pipelines
- Generate pipelines via SDK
- Utilize workflow constructs: loops, branches, conditional execution, variables, parameters, ...



DATA FLOW



- Code-free data transformations that execute in Spark
- Scale-out with Azure Integration Runtimes
- Generate data flows via SDK
- Designers for data engineers and data analysts

The activities in a pipeline define actions to perform on your data.

Pipelines are created code-free using the new V2 Datafactory interface meaning it is easier than ever to create pipelines for Data move and transform activity.

Behind the scenes the managed Apache Spark™ service takes care of the code generation and maintenance.

INGEST



- Multi-cloud and on-prem hybrid copy data
- 90+ native connectors
- Serverless and auto-scale
- Use wizard for quick copy jobs



We will now demo some basic Data Factory functionality in a test environment.



DEMO

Thank You!

Q&A

