# STREAMLINING DATABRICKS

## CI/CD your Notebooks with DevOps Pipelines and orchestrate via Azure Data Factory

**SVENCHIO**

**TechTacoFriday**

Technology spiced with a taco flavor mix

"A discussion on whether use Azure Data Factory and Databricks together is like discussing whether to use a hammer and a screwdriver for any purpose: they are simply different tools, each with unique capabilities and you will achieve a superior result by using the strengths of both. So, let's take advantage of all orchestrating capabilities from Azure Data Factory to support any big data processing, analytics, and machine learning workloads from Databricks!!"

Hector Sven, Linkedin post March 2024

# Workshop's Agenda

**Provision your Workspace infrastructure with BICEP**

**Manage Users with Azure's Entra ID**

**Deploying Notebooks across environments with DevOps YAML Pipelines**

**Orchestrate your Notebooks via Azure Data Factory**

## Hector Sven

✓ From Mexico, living in Norway since Sept 2018

✓ Manager, Data Engineering @Avanade

✓ 20+ years in IT industry

✓ 7 associate certifications on Azure data-related technologies and a DevOps expert

✓ My blog >>>>>>>>>>>>>>>>>>

✓ Contact me
- ❑ hector.lopez@avanade.com
- ❑ svenchio@techtacofriday.com
- ❑ www.linkedin.com/in/svenchio/

## About me…

SVENCHIO

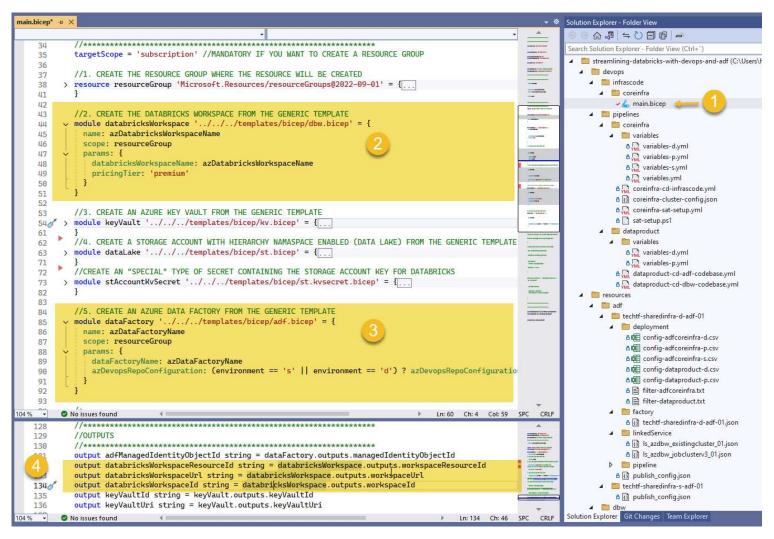TechTacoFriday

Technology spiced with a taco flavor mix

# PROVISIONING DATABRICKS WORKSPACE INFRASTRUCTURE WITH BICEP

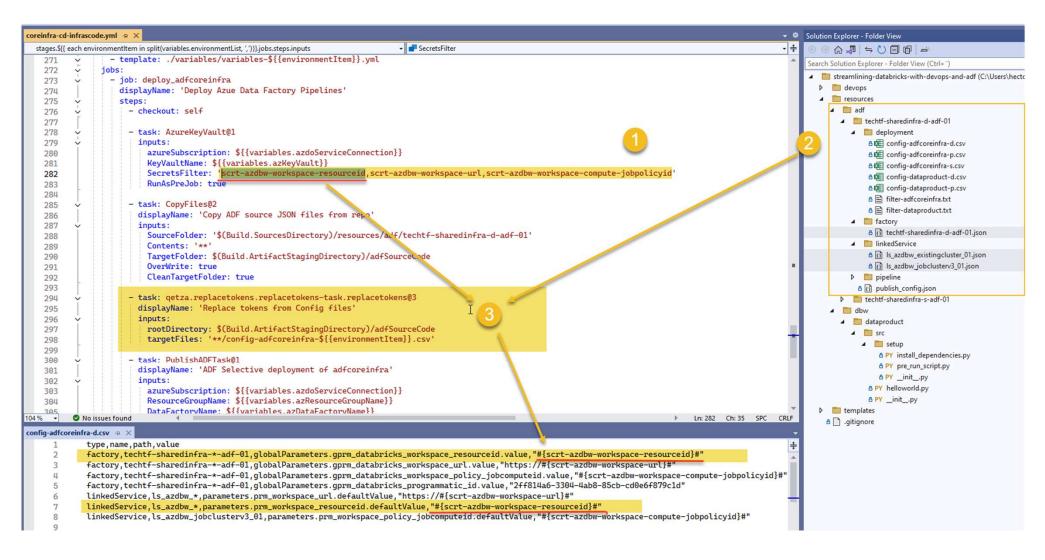**TechTacoFriday**
Technology spiced with a taco flavor mix
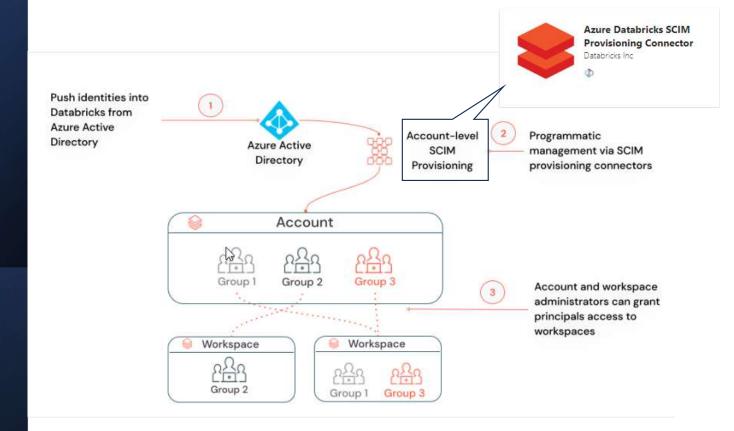
# A few pointers on BICEP

# How it works in a nutshell

**Manage Databricks Users with Azure's Entra ID.**

TechTacoFriday

Technology spiced with a taco flavor mix

# Configuring SCIM for the Account

**Step 1. Check the Requirements**
- ✓ Cloud Application Administrator role in Microsoft Entra ID
- ✓ You must be an Azure Databricks account admin
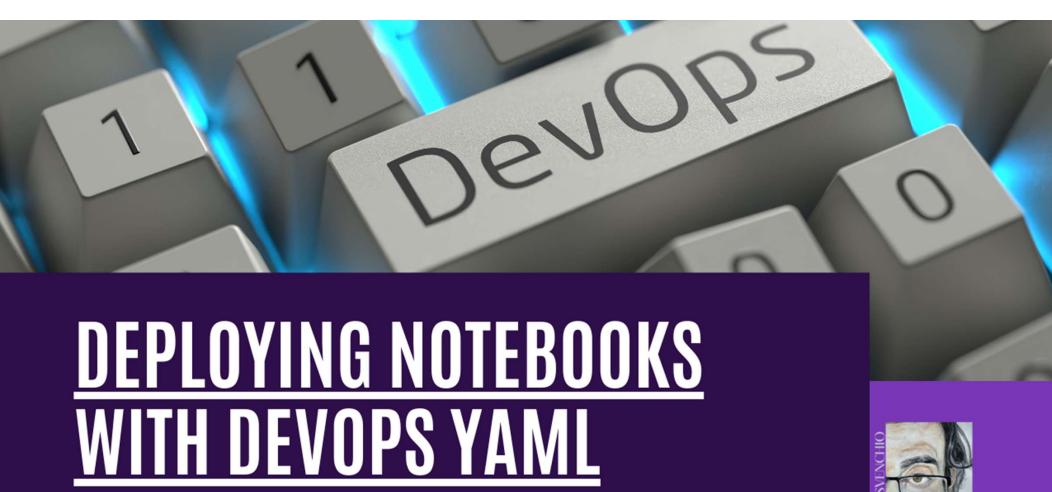- ✓ Azure Databricks account must have the Premium plan

**Step 2. Configure Azure Databricks Account**
- ✓ Enable User Provisioning
- ✓ Copy Account SCIM URL & Token

**Step 3. Create and configure SCIM application**
- ✓ Create the SCIM Enterprise App
- ✓ Configure with Account's SCIM URL & Token

# Quick poll ... which one is your approach?

## Method 1. Databricks Repos with Automated Sync

Repos (a feature in Databricks) are used to sync notebooks from a Git repository directly into a workspace.

## Method 2. CI/CD Pipeline-Based Deployment

This approach involves integrating Databricks into a Continuous Integration/Continuous Deployment (CI/CD) pipeline, often using tools like Jenkins, GitHub Actions, Azure DevOps, or GitLab CI/CD. In this approach, the CI/CD pipeline orchestrates the promotion of notebooks (or other Databricks artifacts) from one environment to another.

## Method 3. API-Based or CLI-Based Deployment

In this approach, deployments are managed via scripts that interact with Databricks through its REST API or CLI. This is similar to the CI/CD approach but can be more lightweight, focusing on custom scripts to promote changes between environments.

## Method 4. Manual approach (aka Copy & Paste)

Copy & paste the content of my workbook from one environment to the other

ORCHESTRATE YOUR
NOTEBOOKS WITH AZURE.

SVENCHIO

**TechTacoFriday**
Technology spiced with a taco flavor mix

# Databricks Linked service & cluster selection

## Job cluster

- ✓ Modular approach
- ✓ Job cluster can be reused in multiple executions
- ✓ Job cluster is terminated when the last chain execution is finished

## Existing cluster

- ✓ Better for schedule tasks that happen regularly
- ✓ Cluster configuration can be stored and reused as well as libraries inhalations
- ✓ Time to live can be configured in the cluster

## Existing instance pool

- ✓ Used when chained executions are needed
- ✓ Same pool can be used in different pipelines (be mindful on the concurrency limits)
- ✓ Lower start up times
- ✓ Job cluster can be configured to use the pool

# Key Advantages of Azure Data Factory (ADF) vs. Databricks Orchestration

| | ADF | Databricks |
|---|---|---|
| Broader Integration Ecosystem | **100+ cloud and on-premise data sources** (e.g., SQL Server, Oracle, SAP, Salesforce, Cosmos DB) | Primarily optimized for running Databricks notebooks and tasks within the Databricks ecosystem. |
| Data Movement & Hybrid Scenarios | Built-in support for **hybrid scenarios** (data migration from on-premises via Integration Runtime and gateways). | Less efficient and harder to manage when orchestrating large-scale data migration from on-premises sources. |
| Low-Code/No-Code Visual Pipeline Development | Visual, drag-and-drop, code-free authoring environment. | Requires notebook-based scripting (Python, Scala, SQL) or JSON-based task configuration |
| Separation of Concerns and Flexibility | Easier to **swap or augment compute technologies** in the future without major re-architecture | Replacing or augmenting compute platforms later can require significant rework. |
| **When to Use Each: Quick Reference** | ✓ Complex workflows involving various Azure and third-party services<br><br>✓ Hybrid (cloud/on-premises) data movement scenarios | ✓ Simple scheduling/execution of Databricks notebooks<br><br>✓ Pure Databricks ecosystems with minimal external interactions |

**Streamlining Databricks: CI/CD your Notebooks with DevOps Pipelines and orchestrate via Azure Data Factory (Series)**

On this series I'm going to show you how to provision your Databricks infrastructure with BICEP and to connect your workspace to Azure's Entra ID to manage users & groups. Furthermore, I'll show you how to deploy your notebooks across Environments with yaml pipelines and orchestrate with ADF

# Download the code here …

🔗 https://www.techtacofriday.com/orchestrate-your-notebooks-via-azure-data-factory/

# Read the article series here

🔗 https://www.techtacofriday.com/streamlining-databricks-cicd-with-bicep-devops-and-adf/

A big shoutout to…

# twoday

# Executions and nice tricks

✓ Execute Azure databricks jobs in sequence or with conditions
  - Use the combination activities and Azure Databricks
    - Cluster reuse
    - Parametrization
    - Rerun
    - Delta live tables
    - Use pipelines as templates
    - Reuse code, triggers and infrastructure

✓ Pools and job clusters
  - Provide workload isolation
  - Reduces pricing
  - Auto termination
  - Faster job cluster creation