

[Get started](#)[Open in app](#)

Beda Tse

86 Followers [About](#)[Follow](#)

Azure DevOps CI/CD with Azure Databricks and Data Factory— Part 1



Beda Tse Feb 28, 2019 · 13 min read

Let's cut long story short, we don't want to add any unnecessary introduction that you will skip anyway.

For whatever reason, you are using Databricks on Azure, or considering using it. Google, or your favourite search engine, has brought you here because you want to explore Continuous Integration and Continuous Delivery using Azure DevOps.

May be it is just me, I always find outdated tutorials when I wanted to learn something new, and by the time you are reading this, this has already outdated as well. But hopefully this series will give you some insight on setting up CI/CD with Azure Databricks.

Without further due, let us begin.

Prerequisites

You need to have an *Azure account*, an *Azure DevOps organisation*, you can leverage either *GitHub* as repository or *Azure Repos* as repository. In this series, we will assume you are using Azure Repos.

You will need to *create a project* on Azure DevOps, together with a *repository*. A sample repository can be found [here](#).

You will need a *git client*, or command line git. We will use command line git throughout the series, thus assuming that you also have a terminal, such as Terminal on Mac, or Git-Bash on Windows.

You will need a *text editor* other than the normal Databricks notebook editor. Visual Studio Code is a good candidate for that. If the text editor have built-in git support, that will be even better.

Checklist

1. Azure Account
2. Azure DevOps Organisation and DevOps Project
3. Azure Service Connections set up inside Azure DevOps Project
4. Git Repos (Assuming you are using Azure Repository)
5. Git Client (Assuming you are using command line Git)
6. Text Editor (e.g. Visual Studio Code, Sublime Text, Atom)

Step 0: Set up your repository and clone it onto your local workstation

0–1. In your Azure DevOps, set up your SSH key or Personal Access Token.

0–2. Create git repo on Azure DevOps with initial README

Create a new repository X

Type ▼

Git

Repository name *

databricks-example

Add a README to describe your repository

Add a .gitignore:

None

Create**Cancel**

0–2. Create your git repo on Azure DevOps inside a project with initial README.

0–3. Locate your git repository URL for git clone

The screenshot shows a repository named 'databricks-example' on Azure DevOps. The repository has one file, 'README.md'. On the right side, there is a 'Clone repository' section. It provides two cloning options: 'HTTPS' (selected) and 'SSH'. The SSH URL is displayed as 'git@ssh.dev.azure.com:v3/bedatse/azure-dataao...'. Below this, there are links for 'Manage SSH keys' and 'Learn more about SSH'. At the bottom of the 'Clone repository' section, there is a note: 'Having problems authenticating in Git? Be sure to get the latest version of [Git for Windows](#) or our plugins for [IntelliJ](#), [Eclipse](#), [Android Studio](#) or [macOS & Linux terminal](#)'.

0–3. Locate your git URL for cloning

0–4. Clone the repository via git using the following command

```
$ git clone <repository-url>
```

```
2. bedatse@Bedas-Macbook: ~/devops-blogs/databricks-example (zsh)
bedatse@Bedas-Macbook ~/devops-blogs git clone git@ssh.dev.azure.com:v3/bedatse/azure-dataops/databricks-example
Cloning into 'databricks-example'...
remote: Azure Repos
remote: Found 3 objects to send. (47 ms)
Receiving objects: 100% (3/3), done.
bedatse@Bedas-Macbook ~/devops-blogs cd databricks-example
bedatse@Bedas-Macbook ~/devops-blogs/databricks-example master ls -al
total 8
drwxr-xr-x  4 bedatse  staff  128 Feb  8 15:07 .
drwxr-xr-x  3 bedatse  staff   96 Feb  8 15:07 ..
drwxr-xr-x 14 bedatse  staff  448 Feb  8 15:15 .git
-rw-r--r--  1 bedatse  staff  936 Feb  8 15:07 README.md
bedatse@Bedas-Macbook ~/devops-blogs/databricks-example master
```

0-4. Now you have cloned your repository locally.

Step 1: Provisioning Azure Databricks and Azure Key Vault with Azure Resource Manager Template

We want to automate the service provisioning or service updates. When you need to set up another set of Databricks, or update anything, you can just update configuration json in the repository, or variables stored on Azure DevOps Pipeline, which we will cover in the next steps. Azure DevOps Pipeline will take care of everything else for you.

1-1. Copy `template.json` `parameters.json` `azure-pipeline.yml` `notebook-run.json.tpl` from [this commit](#) of the example repository, put them into your repository local folder.

1-2. Stage the changed file in git, commit and push it onto the Azure Repo.

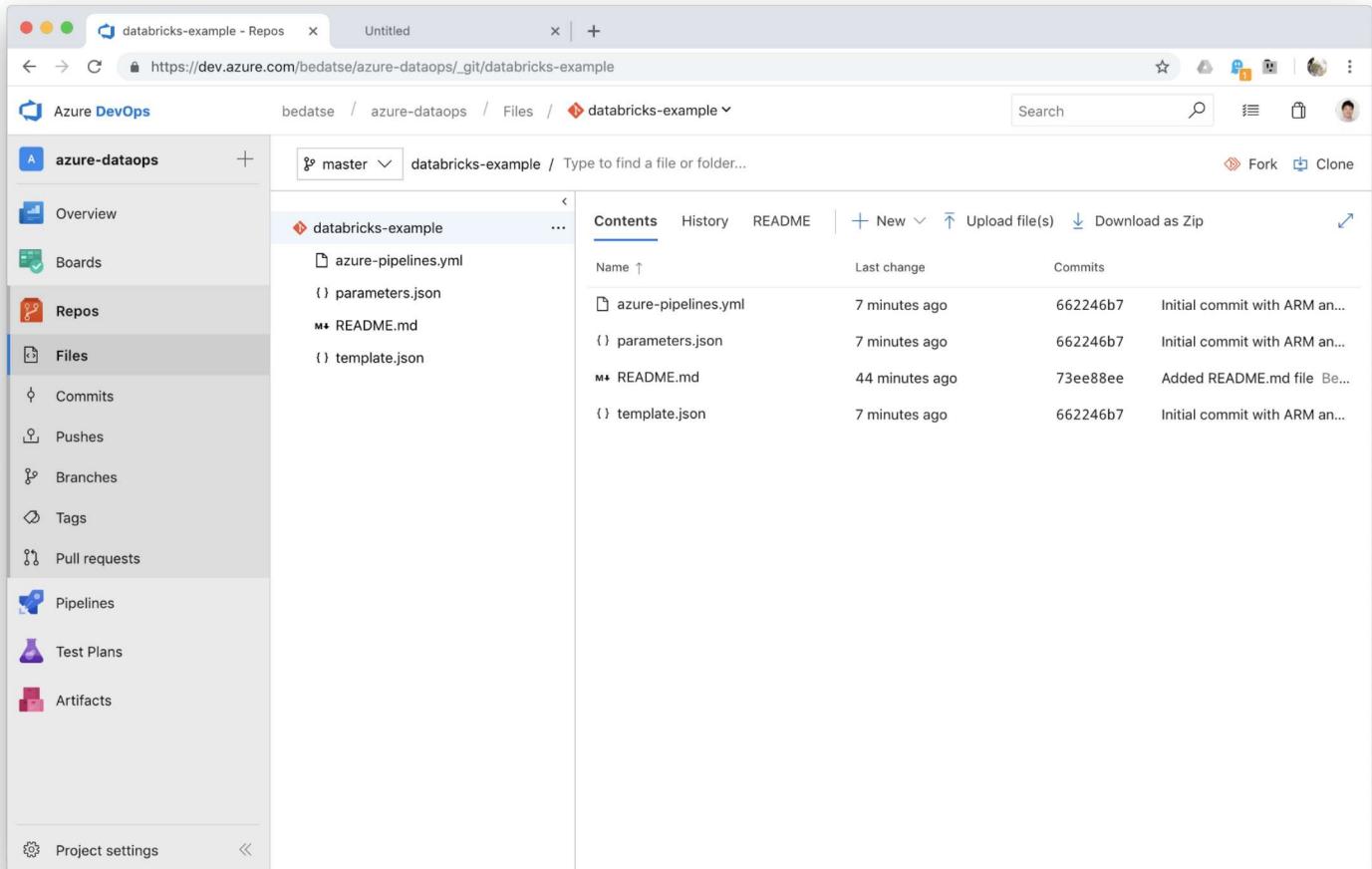
```
$ git add -A  
$ git commit -m '<your-commit-message>'  
$ git push
```

The screenshot shows a terminal window with the following session:

```
2. bedatse@Bedas-Macbook: ~/devops-blogs/databricks-example (zsh)  
x ..ricks-example (zsh) #1  
      template.json  
  
nothing added to commit but untracked files present (use "git add" to track)  
bedatse@Bedas-Macbook ~ ~/devops-blogs/databricks-example ↵ master ➔ git add -A  
bedatse@Bedas-Macbook ~ ~/devops-blogs/databricks-example ↵ master + ➔ git commit -m 'Initial commit with ARM and build pipeline'  
[master 662246b] Initial commit with ARM and build pipeline  
 3 files changed, 196 insertions(+)  
  create mode 100644 azure-pipelines.yml  
  create mode 100755 parameters.json  
  create mode 100755 template.json  
bedatse@Bedas-Macbook ~ ~/devops-blogs/databricks-example ↵ master ➔ git push  
  
Counting objects: 5, done.  
Delta compression using up to 8 threads.  
Compressing objects: 100% (5/5), done.  
Writing objects: 100% (5/5), 1.94 KiB | 1.94 MiB/s, done.  
Total 5 (delta 0), reused 0 (delta 0)  
remote: Analyzing objects... (5/5) (15 ms)  
remote: Storing packfile... done (76 ms)  
remote: Storing index... done (25 ms)  
To ssh.dev.azure.com:v3/bedatse/azure-dataops/databricks-example  
 73ee88e..662246b master -> master  
bedatse@Bedas-Macbook ~ ~/devops-blogs/databricks-example ↵ master
```

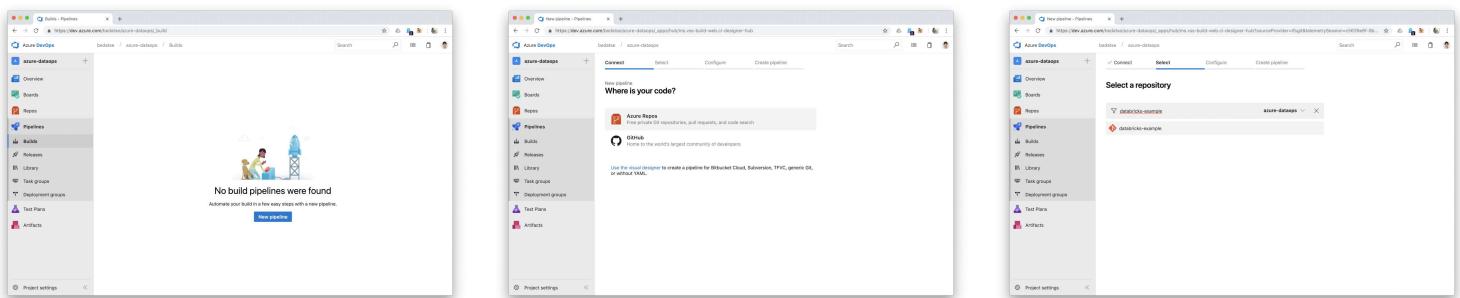
The terminal uses color-coded arrows to indicate the flow of commands: blue for navigating to the directory (~), green for switching branches (master), yellow for committing, and red for pushing.

1–2. Commit and Push infrastructure code and build pipeline code onto repository



1–2. After pushing code back into repository, it should look like this.

1–3. Create your build pipeline, go to **Pipelines > Builds** on the sidebar, click **New Pipeline** and select **Azure DevOps Repo**. Select your repository and review the pipeline `azure-pipeline.yml` which has already been uploaded in step 1–2. Click **Run** to run the build pipeline for the first time.



1–3–1 Create new Build Pipeline

The left screenshot shows the 'Create pipeline' step in the Azure DevOps UI. It displays the YAML code for 'azure-pipelines.yml':

```

1 trigger:
2   - master
3
4 pool:
5   - vmImage: 'Ubuntu-16.04'
6
7 stages:
8   - build: parameters.json template.json "$Build.ArtifactStagingDirectory"
9     displayName: 'Include Azure Resource Manager templates into Build Artifacts'
10    tasks:
11      - task: PublishBuildArtifacts@1
12        inputs:
13          pathToPublish: '$Build.ArtifactStagingDirectory'
14          artifactName: arm_template
15

```

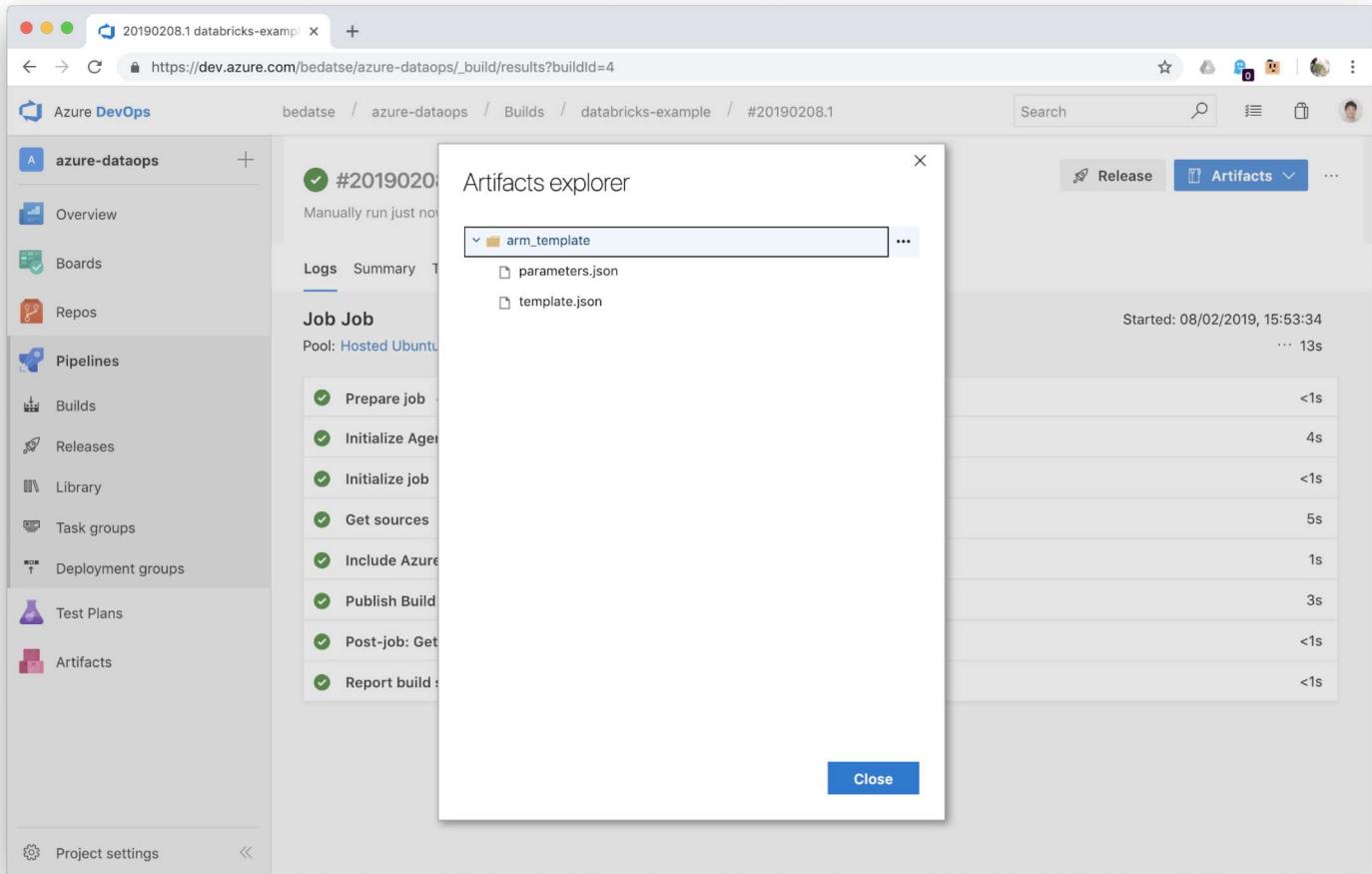
The right screenshot shows the execution summary of a build pipeline named '#20190208.1: Initial commit with ARM and build pipeline'. The summary table includes columns for Step, Status, and Duration.

| Step | Status | Duration |
|---|-----------|----------|
| Prepare job | succeeded | <1s |
| Initialize Agent | succeeded | 4s |
| Initialize job | succeeded | <1s |
| Get sources | succeeded | 5s |
| Include Azure Resource Manager templates into Build Artifacts | succeeded | 1s |
| Publish Build Artifacts | succeeded | 3s |
| Post-job: Get sources | succeeded | <1s |
| Report build status | succeeded | <1s |

1–3–2 Review the content of the pipelines and execution result

The build pipeline currently only do one thing, which is to pack the Azure Resource Manager JSONs into a build artifact, which can be consumed on later steps for deployment. Let take a look what is inside of the artefact now.

In your build, click **Artifacts > arm_templates**. The details of the artifact will be displayed.



The artifact **arm_template** contains ARM JSON files.

1–4. Create variable group for your deployment. You don't want to hardcode your variables inside the pipeline, such that you can make it reusable in another project or environment with least effort. First, let's create a group for your Project, storing all variables that would be the same across all environments.

Go to **Pipelines > Library**, Click on **+Variable group**. Type in your variable group name, as an example, we are using *Databricks Pipeline* as the variable group name. Add a variable with *Name* `project_name` with *Value* `databricks-pipeline` . Save the changes after you are done.

The screenshot shows the Azure DevOps Library - Pipelines page. The left sidebar is titled 'azure-dataops' and includes links for Overview, Boards, Repos, Pipelines, Builds, Releases, Library (which is selected), Task groups, Deployment groups, Test Plans, and Artifacts. The main content area is titled 'Databricks Pipeline' under 'Variable group'. It has tabs for Variable group, Save, Clone, Security, and Help. Under 'Properties', there is a 'Variable group name' field containing 'Databricks Pipeline' and a 'Description' field which is empty. There are two toggle switches: 'Allow access to all pipelines' (on) and 'Link secrets from an Azure key vault as variables' (off). Under 'Variables', there is a table with one entry: 'project_name' with a value of 'databricks-pipeline'. A '+ Add' button is available for adding more variables.

1–4–1 Project-wide Variable Group

Create another group named *Dev Environment Variables*, this one will have more variables in it, as listed below.

- `databricks_location` : <databricks location>
- `databricks_name` : <databricks name>
- `deploy_env` : <deployment environment name>
- `keyvault_owner_id` : <your user object ID in Azure AD>
- `keyvault_name` : <key vault name for storing databricks token>
- `rg_groupname` : <resource group name>
- `rg_location` : <resource group location>

- tenant_id :<your Azure AD tenant ID>

The screenshot shows the Azure DevOps Library - Pipelines page. On the left, there's a sidebar with options like Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The 'Library' option is selected. In the main area, it says 'Library > Dev Environment Variables'. There are tabs for Variable group (selected), Save, Clone, Security, and Help. Below that is a 'Properties' section with a 'Variable group name' set to 'Dev Environment Variables' and a 'Description' field. There are two toggle buttons: 'Allow access to all pipelines' (which is turned on) and 'Link secrets from an Azure key vault as variables' (which is turned off). Below this is a 'Variables' section with a table:

| Name ↑ | Value |
|---------------------|------------------------|
| databricks_location | southeastasia |
| databricks_name | sample-databricks-dev |
| deploy_env | dev |
| keyvault_name | bpdo-databricks-kv-dev |
| keyvault_owner_id | 359c [REDACTED] 815e |
| rg_groupname | bp-dataops-dev |
| rg_location | southeastasia |
| tenant_id | 8286 [REDACTED] 5c5d |

At the bottom of the table is a '+ Add' button.

1–5. Create a new release pipeline, Go to **Pipelines** > **Releases**, click on **+ New**. Select start with an ***Empty job*** when you are asked to select a template. Name your stage ***Dev Environment***, in future tutorials, we will be cloning this stage for ***Staging Environment*** and ***Production Environment***.

The three screenshots illustrate the steps to create a new release pipeline:

- The first screenshot shows the 'Releases' page in Azure DevOps. It has a sidebar with 'Overview', 'Boards', 'Repos', 'Pipelines', 'Builds', 'Releases' (selected), 'Task groups', 'Deployment groups', 'Test Plans', and 'Artifacts'. The main area says 'No release pipeline selected' and 'Select a release pipeline to view its releases'.
- The second screenshot shows the 'Select a template' dialog. It lists several options under 'Featured':
 - Azure App Service deployment
 - Deploy a Node.js app to Azure App Service
 - Deploy a Node.js app to Azure Web App
 - Deploy a PHP app to Azure App Service and App Service Environment
 - Deploy a Python app to Azure App Service and App Service Environment
 - Deploy a Python app to Azure App Service
 - Deploy a Python app to Azure App Service and App Service Environment
 - Deploy to a Kubernetes Cluster
 - ES vehicle and SQL database deployment
- The third screenshot shows the 'New release pipeline' configuration page. It has a sidebar with 'Overview', 'Boards', 'Repos', 'Pipelines' (selected), 'Builds', 'Releases' (selected), 'Task groups', 'Deployment groups', 'Test Plans', and 'Artifacts'. The main area shows a single stage named 'Dev Environment' with a status of '1 job, 0 task'. A 'Stage definition' dialog is open, showing the stage name 'Dev Environment' and the stage owner 'Beda Tse'.

1–5–1 Create New Release Pipeline

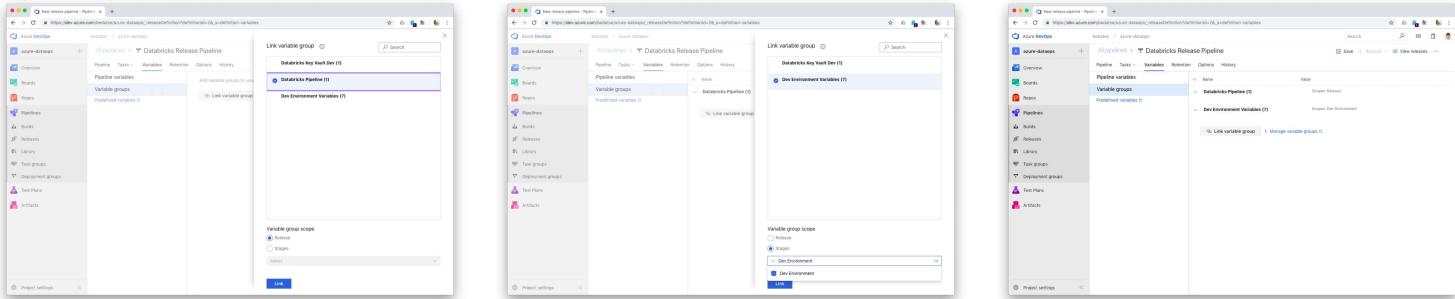
Add your build artifact from your repository as the source artifact. In this example, we will add from databricks example. Click **+Add** next to *Artifacts*.

Select **Build** as *Source type*, select your *Project* and *Source*. Select **Latest** as *Default version*, keep *Source alias* unchanged. Click **Add** when you are done.

1–5–4 Add an artifact to the release pipeline as trigger.

Before moving onto specifying *Tasks* in the *Dev Environment Stage*, let's link the variable group with the release pipeline and the Dev Environment Stage.

Go to **Variables > Variable groups**, Click **Link variable group**. Link your *Databricks Pipeline* created in step 1–4 with scope set to *Release*. Link your *Dev Environment Variables* created in step 1–4 with scope set to *Stages*, and apply to *Dev Environment stage*.



1–5–5 Link Databricks Pipeline Project Variable Group with Release scope

With variable groups linked, we are ready for setting up tasks in the ***Dev Environment*** stage. Click **Tasks > Agent job**, review the settings in there.

1–5–6 Empty Agent Job

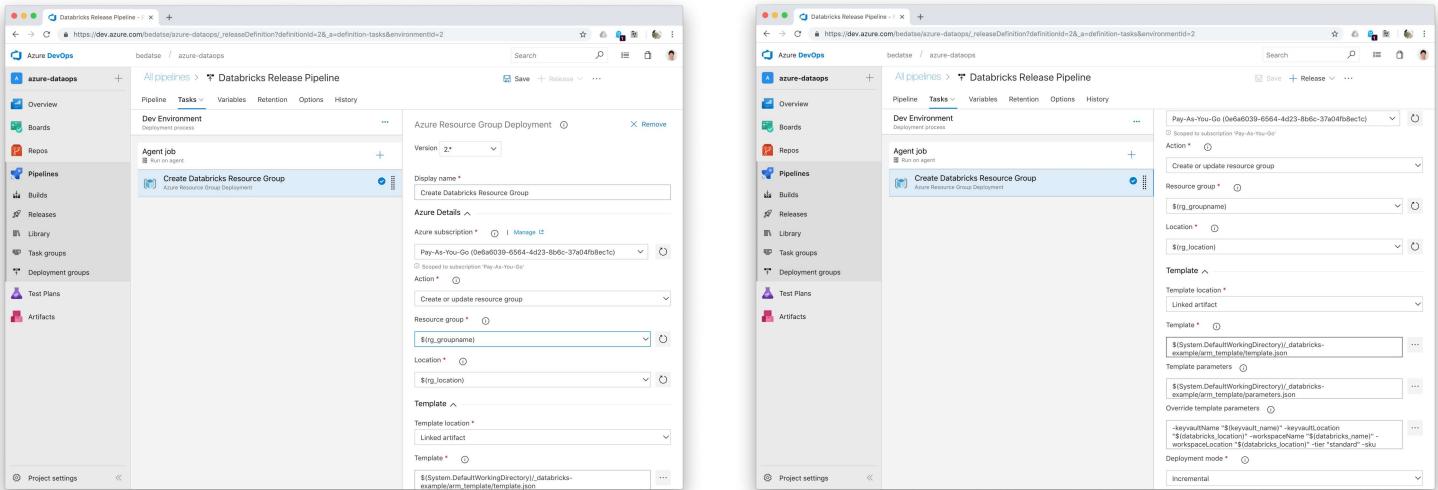
Click the + sign next to the Agent job, add an **Azure Resource Group Deployment**

task.

1–5–7 Add task to Deployment stage

Now, we can configure **Azure Resource Group Deployment** task. Select the service connections from previous step, keep **Action** as **Create or update resource group**. For **Resource group name** and **location**, type in `$(rg_groupname)` and `$(rg_location)` respectively.

For **Template** and **Template Parameters**, click on the **More** action button next to the text field, select **template.json** and **parameters.json** inside `_databricks-example/arm_template`. Before we set our override template parameters, let us set the **Deployment mode** to **Incremental**.



1–5–8 Configuring Resource Group Task

The most challenging part in this section is ***Override template parameters***, we have made this simple for you. Just copy the following snippet into the text field for now. This will allow you to override the default value specified in the *parameters file* by the value specified in *Variable Group*.

```
-keyvaultName "$(keyvault_name)" -keyvaultLocation
"$(databricks_location)" -workspaceName "$(databricks_name)" -
workspaceLocation "$(databricks_location)" -tier "standard" -sku
"Standard" -tenant "$(tenant_id)" -enabledForDeployment false -
enabledForTemplateDeployment true -enabledForDiskEncryption false -
networkAcls
{"defaultAction":"Allow", "bypass":"AzureServices", "virtualNetworkRule
s": [], "ipRules": []}
```

After all these, save your release pipeline and we are ready to create a release.

1–6. Create a release by going back to **Pipelines > Releases** screen. Click on **Create a release** button, then click Create. Your release will then be queued.

Create a new release
Databricks Release Pipeline

Pipeline ↗
Click on a stage to change its trigger from automated to manual.

Dev Environment

Stages for a trigger change from automated to manual. ⓘ

Artifacts ↗
Select the version for the artifact sources for this release

| Source alias | Version |
|---------------------|------------|
| _databricks-example | 20190210.1 |

Release description

Create **Cancel**

1–6–1 Create a release

Releases - Pipelines

Databricks Release Pipeline - Release-1

Release Stages

Manually triggered 2019-08-29 10:00:00 UTC

Dev Environment

Agent job

Initiate Agent - succeeded

Initiate job - succeeded

Download Artifacts - succeeded

Create Databricks Resource Group

Waiting for console output from an agent...

Release - 1

Deployment process

Agent job

Initiate Agent - succeeded

Initiate job - succeeded

Download Artifacts - succeeded

Create Databricks Resource Group

Waiting for console output from an agent...

Release - 1

Deployment process

Agent job

Initiate Agent - succeeded

Initiate job - succeeded

Download Artifacts - succeeded

Create Databricks Resource Group

Waiting for console output from an agent...

Release - 1

Deployment process

Agent job

Initiate Agent - succeeded

Initiate job - succeeded

Download Artifacts - succeeded

Create Databricks Resource Group

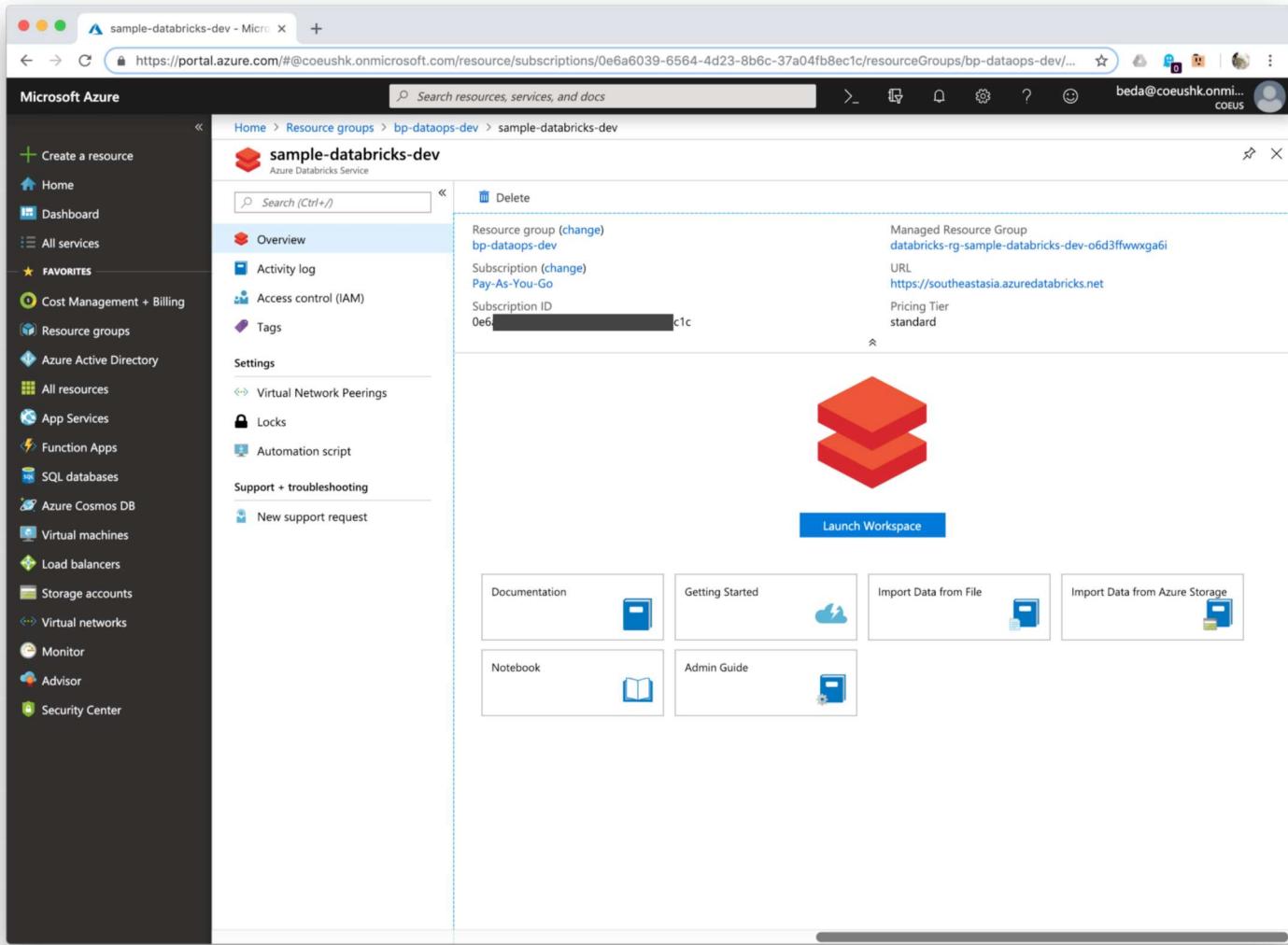
Waiting for console output from an agent...

Release - 1

1–6–2 Wait for provisioning result

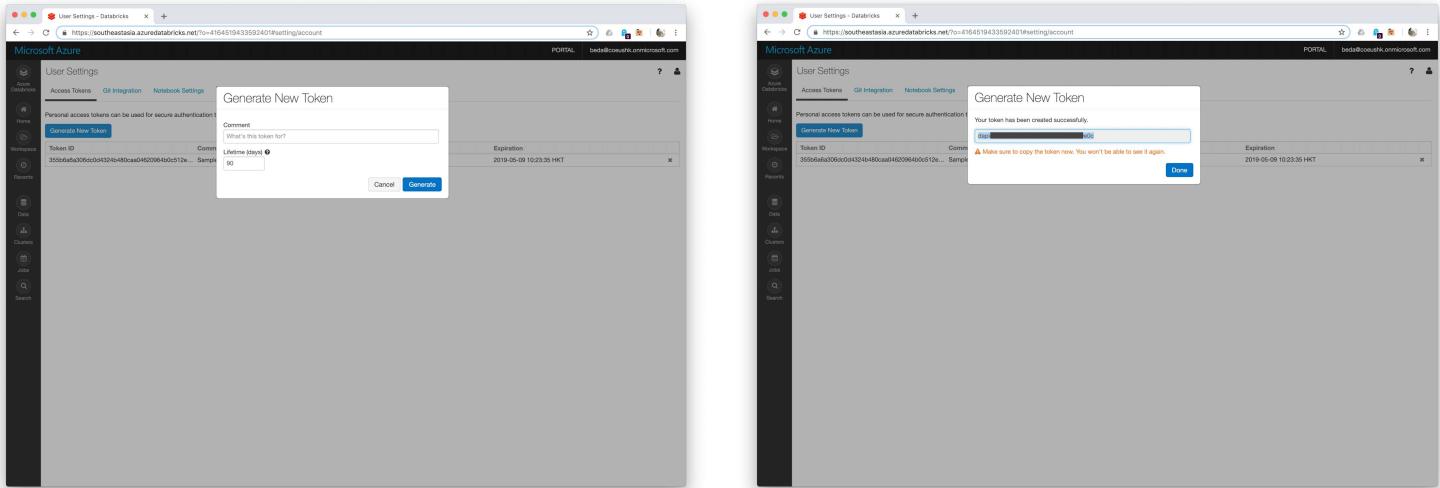
Step 2: Generate Azure Databricks API Token and store the token into Azure Key Vault

2–1. Access Azure Portal, look for the newly created resource group and Databricks, and launch Databricks Workspace as usual.



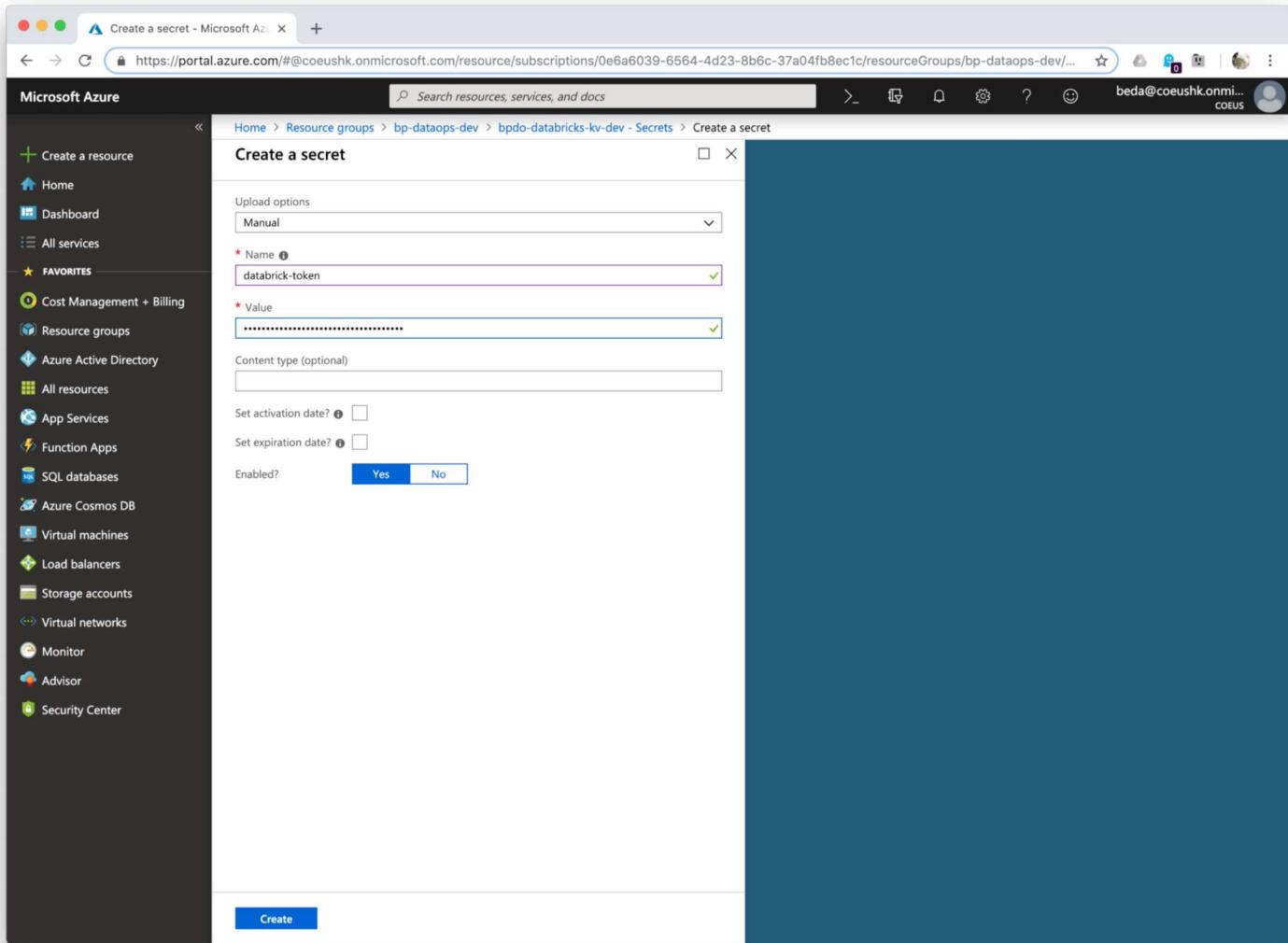
2–1 Access Databricks Workspace via Azure Portal

2–2. After logging into the workspace, click the **user icon** on the top right corner, select **User Settings**. Click **Generate New Token**, give it a meaningful comment and click **Generate**. We will use this token in our pipeline for Notebook deployment. Your token will only be displayed once, make sure you do not close the dialog or browser before you have copied it into key vault.



2-2 Generate new token

2-3. In another browser window, open Azure Portal, navigate to Azure Key Vault under the newly created Resource group. Access the **Secrets** tab, click **Generate/Import**. Set the **Name** as `databricks-token`, and copy the newly generated token into **Value**. Click **Create** to save the token inside Azure Key Vault securely. Now you can safely close the



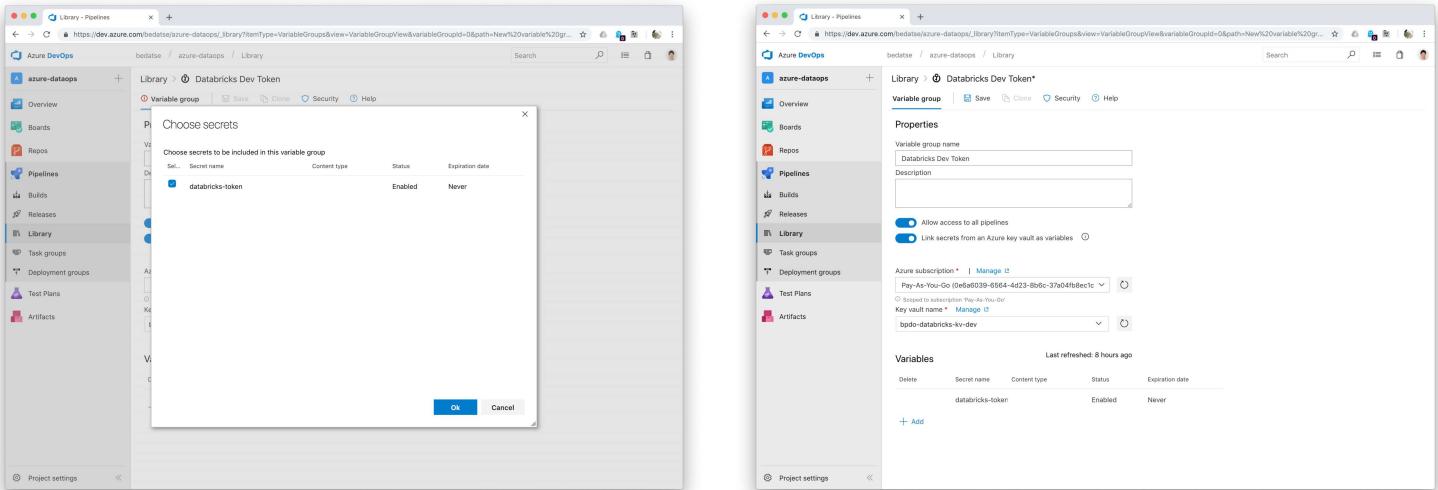
2–3 Save Databricks token into Azure Key Vault

2–4. While we are in the Databricks workspace, also go to **Git Integration** tab and check **Git provider** setting, make sure it is set to **Azure DevOps Services**, or to the repository of your choice.

The screenshot shows the 'User Settings' page for Azure Databricks. The left sidebar includes icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main content area is titled 'User Settings' and has tabs for 'Access Tokens', 'Git Integration' (which is selected), and 'Notebook Settings'. Under 'Git Integration', there is a section titled 'Connect to your Git Account' with instructions for GitHub and Bitbucket. A dropdown menu labeled 'Git provider' is set to 'Azure DevOps Services'. A 'Change settings' button is at the bottom.

2–4 Ensure git integration settings is set to Azure DevOps Services

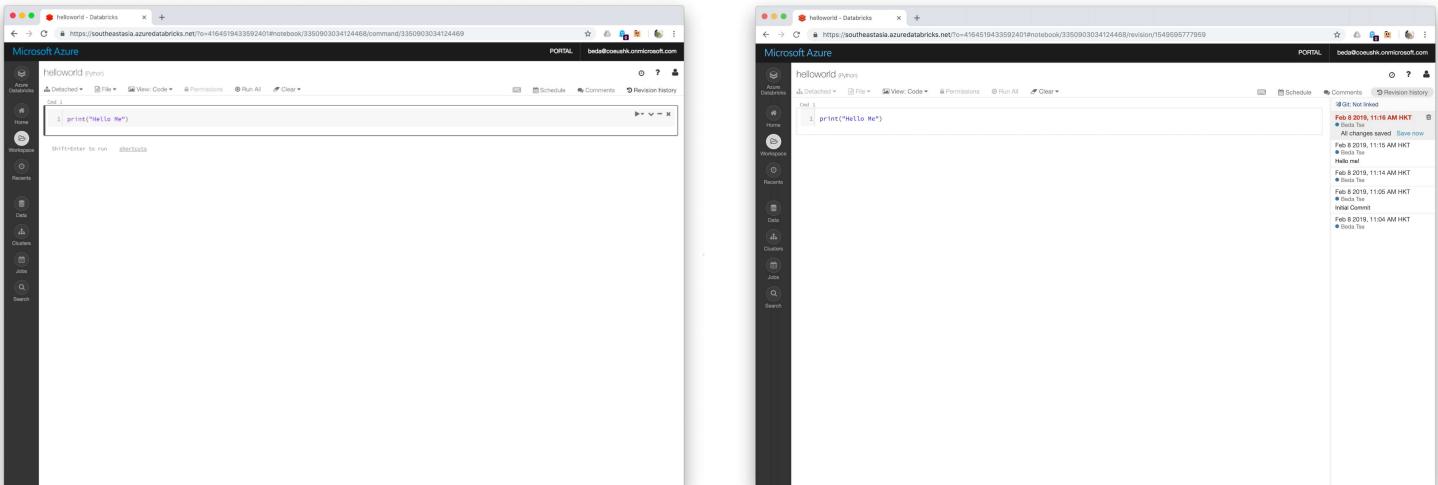
2–5. Go to Azure DevOps Portal, go to **Pipelines > Library**. Click **+Variable Group** to create new Variable Group. This time we are linking an Azure Key Vault into Azure DevOps as variable group. This allows Azure DevOps to obtain token from Azure Key Vault securely for deployment. Name the variable group as **Databricks Dev Token**, select **Link secrets from an Azure key vault as variables**. Select the correct Azure subscription service connections and Key vault respectively. Click **+Add** and select **databricks-token** in the **Choose secrets** dialog. Click **Save**.



Step 3: Link your workbook development with Source code repository

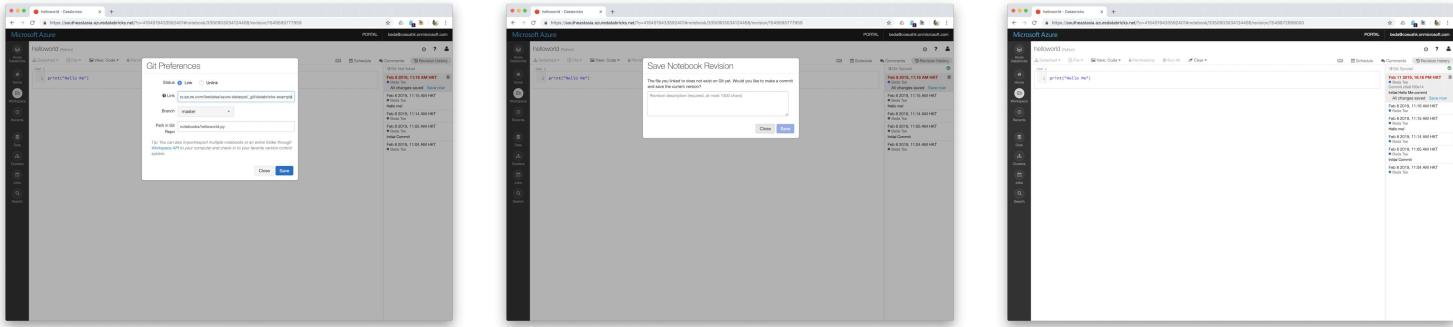
Databricks workspaces integrate with git seamlessly as an IDE. We have in previous steps 2–4 set up the integration between Databricks and a source code repository. We can now link the workbook with the repository and commit into a repository directly.

3–1. Open your notebook as usual, notice **Revision history** on the right top section of the screen. Click on **Revision history** to bring up the version history side panel.



3–1 Version history side panel.

3–2. Click on **Git: Not Linked** to update **Git Preferences**. Link your workbook to *Azure DevOps Repo*, which should be the URL of your git repository, and set the **Path in Git Repo** to the location which you want Databricks to save your notebook inside the repository.



3–2 Link workbook with repository

Azure DevOps GitHub interface showing a committed change to a Databricks notebook. The commit message is 'Initial Hello Me commit' by 'Beda Tse'.

3–2 Committed change pushed into git repo.

3–3. The committed change is pushed into git repository. What does that mean? That means it will trigger the build pipeline. With a little bit of further configuration, we can actually update the build pipeline to package this notebook into a deployable package, and use it to trigger a deployment pipeline. Now download the *azure-pipelines.yml*

from [this commit](#), replace the original `azure-pipelines.yml` from step 1–1, commit and push the change back to the repository.

The screenshot shows the Azure DevOps interface for a pipeline named 'azure-dataops'. A specific build, '#20190211.2: add notebook build handling', is displayed. The build was triggered today at 23:28 for the user 'Beda Tse' on the 'databricks-example' branch. The build status is 'succeeded'. The pipeline log shows the following tasks:

- Prepare job - succeeded (duration <1s)
- Initialize Agent - succeeded (duration 5s)
- Initialize job - succeeded (duration <1s)
- Get sources - succeeded (duration 6s)
- Include Azure Resource Manager templates into Build Artifacts - succeeded (duration 1s)
- Prepare Notebook Build Artifacts - succeeded (duration <1s)
- Publish ARM Template Build Artifacts - succeeded (duration 3s)
- Publish Notebook Build Artifacts - succeeded (duration 5s)
- Post-job: Get sources - succeeded (duration <1s)
- Report build status - succeeded (duration <1s)

3–3 Build triggered by the pipeline update.

Step 4: Deploy the version controlled Notebook onto Databricks for automated tests

Since we have prepared our notebook build package, let us complete the flow by deploying it onto the Databricks that we have created and execute a run from the pipeline.

4–1. Go to **Pipelines > Library**, edit the project based Variable group **Databricks Pipeline**, we need to specify a variable here such that the release pipeline will pickup from the variable which notebook to deploy. Add `notebook_name` variable with value `helloworld`. Click **Save** after it is done.

The screenshot shows the Azure DevOps Library - Pipelines page. On the left, there's a sidebar with options like Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The Pipelines option is selected. The main area shows a 'Variable group' named 'Databricks Pipeline'. It has a 'Properties' section with fields for Variable group name (set to 'Databricks Pipeline') and Description (empty). There are two toggle switches: 'Allow access to all pipelines' (on) and 'Link secrets from an Azure key vault as variables' (off). Below this is a 'Variables' section with a table:

| Name ↑ | Value |
|---------------|---------------------|
| notebook_name | helloworld |
| project_name | databricks-pipeline |

At the bottom of the table is a '+ Add' button.

4–1 Update Variable group with notebook_name variable

4–2. Go to **Pipelines > Releases**, select **Databricks Release Pipeline**, Click **Edit**. Navigate to the **Tasks** tab. Add **Use Python Version** task and drag it above the original **Create Databricks Resource Group** task. No further configuration is needed.'

4–3. Add Bash Task at the end of the job. Rename it to *Install Tools*. Select **Type** as **Inline**, copy the following scripts to the **Script** text area. This is to install the needed python tools for deploying notebook onto Databricks via command line interface.

```
python -m pip install --upgrade pip setuptools wheel databricks-cli
```

4–3 Configure Install Tools Task

4–4. Add Bash Task at the end of the job. Rename it to *Authenticate with Databricks CLI*. Select **Type** as **Inline**, copy the following scripts to the **Script** text area. The variable `databricks_location` is obtained from variable group defined inside the pipeline, while `databricks-token` is obtained from variable group linked with Azure Key Vault.

```
databricks configure --token <<EOF
https://$(databricks_location).azuredatabricks.net
$(databricks-token)
EOF
```

The screenshot shows the Azure DevOps interface for managing a release pipeline. On the left, there's a sidebar with options like Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The main area displays the 'Databricks Release Pipeline' under 'All pipelines'. The pipeline consists of one 'Agent job' which contains several tasks. One task, 'Authenticate with Databricks CLI', is currently selected. Its configuration includes a 'Display name' of 'Authenticate with Databricks CLI', a 'Type' of 'Inline', and a 'Script' block containing the command `databricks configure --token <<EOF https://\$databricks_location).azuredatabricks.net \$databricks-token EOF`. Other tasks in the job include 'Use Python 3.x', 'Create Databricks Resource Group', 'Install Tools', and another 'Authenticate with Databricks CLI' task.

4–4 Configure CLI Authentication Task

4–5. Add Bash Task at the end of the job. Rename it to *Upload Notebook to Databricks*. Select **Type** as **Inline**, copy the following scripts to the **Script** text area. The variable `notebook_name` is retrieved from the release scoped variable group.

```
databricks workspace.mkdirs /build
databricks workspace.import --language PYTHON --format SOURCE --
overwrite _databricks-
example/notebook/${notebook_name}-${Build.SourceVersion}.py
/build/${notebook_name}-${Build.SourceVersion}.py
```

The screenshot shows the Azure DevOps Pipeline Editor interface. On the left, there's a sidebar with options like Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The main area shows a pipeline named 'Databricks Release Pipeline'. It consists of two steps: 'Dev Environment' (Deployment process) and 'Agent job' (Run on agent). The 'Agent job' step contains several tasks: 'Use Python 3.x' (Use Python Version), 'Create Databricks Resource Group' (Azure Resource Group Deployment), 'Install Tools' (Bash), 'Authenticate with Databricks CLI' (Bash), and 'Upload Notebook to Databricks' (Bash). The 'Upload Notebook to Databricks' task is currently selected. In the details pane on the right, you can see its configuration: it's a 'Bash' task with version 3.*, a display name 'Upload Notebook to Databricks', and a type set to 'Inline'. The 'Script' field contains the following command:

```
#!/bin/bash
databricks workspace mkdirs /build
databricks workspace import --language PYTHON --format SOURCE --overwrite \
    _databricks-example/notebook/${notebook_name}-$(Build.SourceVersion).py \
    /build/${notebook_name}-$(Build.SourceVersion).py
```

Below the script, there are sections for Advanced, Control Options, Environment Variables, and Output Variables.

4–5 Configure Upload Notebook Task

4–6. Add Bash Task at the end of the job. Rename it to *Create Notebook Run JSON*. Select **Type** as **Inline**, copy the following scripts to the **Script** text area. This is to prepare a job execution configuration for the test run, using the template `notebook-run.json.tpl`.

```
# Replace run name and deployment notebook path
cat _databricks-example/notebook/notebook-run.json.tpl | jq
'.run_name = "Test Run - $(Build.SourceVersion)" |
.notebook_task.notebook_path =
"/build/${notebook_name}-$(Build.SourceVersion).py"' >
${notebook_name}-$(Build.SourceVersion).run.json

# Check the Content of the generated execution file
cat ${notebook_name}-$(Build.SourceVersion).run.json
```

The screenshot shows the Azure DevOps Pipelines interface for a 'Databricks Release Pipeline'. The pipeline consists of several tasks:

- Dev Environment
- Agent job
- Use Python 3.x
- Create Databricks Resource Group
- Install Tools
- Authenticate with Databricks CLI
- Upload Notebook to Databricks
- Create Notebook Run JSON** (selected)
- Run Notebook on Databricks
- Wait for Databricks Run to complete

The 'Create Notebook Run JSON' task is currently selected. Its configuration includes:

- Type: Inline
- Script:

```
# Replace run name and deployment notebook path
cat _databricks-example/notebook/notebook-run.json.tpl | jq '.run_name = "Test Run - $(Build.SourceVersion)"', .notebook_task.notebook_path = "build/${notebook_name}-${Build.SourceVersion}.py" > ${notebook_name}-${Build.SourceVersion}.run.json

# Check the Content of the generated execution file
cat ${notebook_name}-${Build.SourceVersion}.run.json
```

4–6 Configure Notebook Run JSON Creation task

4–7. Add Bash Task at the end of the job. Rename it to *Run Notebook on Databricks*.

Select **Type** as **Inline**, copy the following scripts to the **Script** text area. This is to execute the notebook prepared in the Build pipeline, i.e. committed by you thru the Databricks UI, via Job Cluster.

```
echo ##vso[task.setvariable variable=RunId;
isOutput=true;]`databricks runs submit --json-file
${notebook_name}-$(Build.SourceVersion).run.json | jq -r .run_id`"
```

You might have noticed there is weird template here with `##vso[task.setvariable variable=RunId; isOutput=true;]`. This is to save the `run_id` from the output of the

databricks runs submit command into Azure DevOps as variable `RunId`, such that we can reuse that run id in next steps.

The screenshot shows the Azure DevOps interface for a pipeline named 'Databricks Release Pipeline'. The left sidebar is open, showing the 'Pipelines' section. The main area displays the pipeline tasks under the 'Tasks' tab. One task, 'Run Notebook on Databricks', is selected and highlighted with a blue border. This task is a Bash script. The configuration pane on the right shows the following details:

- Type:** Inline
- Script:**

```
echo "##vso[task.setvariable variable=RunId; isOutput=true;]`databricks runs submit --json-file ${notebook_name}-${Build.SourceVersion}.run.json | jq -r .run_id`"
```

4–7 Configure Notebook Execution Task

4–8. Add Bash Task at the end of the job. Rename it to *Wait for Databricks Run to complete*. Select **Type** as **Inline**, copy the following scripts to the **Script** text area. This is to wait for the previously executing Databricks job and get the execution state from the run result.

```
echo "Run Id: $(RunId)"

# Wait until job run finish
while [ "`databricks runs get --run-id $(RunId) | jq -r '.state.life_cycle_state'" != "INTERNAL_ERROR" ] && [ "`databricks runs get --run-id $(RunId) | jq -r '.state.result_state'" == "null" ]
```

```

]
do
echo "Waiting for Databrick job run ${RunId} to complete, sleep for
30 seconds"
sleep 30
done

# Print Run Results
databricks runs get --run-id ${RunId}

# If not success, report failure to Azure DevOps
if [ `databricks runs get --run-id ${RunId} | jq -r
'.state.result_state'` != "SUCCESS" ]
then
echo "#vso[task.complete result=Failed;]Failed"
fi

```

The screenshot shows the Azure DevOps interface for configuring a pipeline. The left sidebar is titled 'azure-dataops' and lists various DevOps services: Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The 'Pipelines' option is selected.

The main area shows the 'Databricks Release Pipeline' under 'All pipelines'. The 'Tasks' tab is selected. The pipeline consists of the following tasks:

- Dev Environment**: Deployment process
- Agent job**: Run on agent
- Use Python 3.x**: Use Python Version
- Create Databricks Resource Group**: Azure Resource Group Deployment
- Install Tools**: Bash
- Authenticate with Databricks CLI**: Bash
- Upload Notebook to Databricks**: Bash
- Create Notebook Run JSON**: Bash
- Run Notebook on Databricks**: Bash
- Wait for Databricks Run to complete**: Bash

The 'Wait for Databricks Run to complete' task is currently selected. Its configuration details are displayed on the right side:

- Type**: Bash
- Version**: 3.*
- Display name**: Wait for Databricks Run to complete
- Type**: Inline
- Script** (content):

```

echo "Run Id: ${RunId}"

# Wait until job run finish
while [ `databricks runs get --run-id ${RunId} | jq -r '.state.life_cycle_state'` != "INTERNAL_ERROR" ] && [ `databricks runs get --run-id ${RunId} | jq -r '.state.result_state'` == "null" ]
do
echo "Waiting for Databrick job run ${RunId} to complete, sleep for 30 seconds"
sleep 30
done

```
- Advanced**, **Control Options**, **Environment Variables**, and **Output Variables** sections are also visible.

4–8 Configure Databricks Run task

4–9. Remember we have added Databricks token from Azure Key vault? It is now to put it in use. Access the **Variables** Tab, click **Variable groups**, Link the variable group **Databricks Dev Token** with **Dev Environment** Stage.

The screenshot shows the Azure DevOps interface for a 'Databricks Release Pipeline'. The left sidebar lists various project sections like Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The main area shows the pipeline navigation path: All pipelines > Databricks Release Pipeline. The 'Variables' tab is selected. Under 'Variable groups', there is a list with one item: 'Databricks Dev Token (1)'. Below this, under 'Variable group scope', 'Stages' is selected, and 'Dev Environment' is chosen from a dropdown menu. A 'Link' button is visible at the bottom of this section.

4–9 Link the Token variable group with Environment

4–10. Save the Release Pipeline, and *create a release* to test the new pipeline.

The screenshot shows the Azure DevOps interface for a 'Databricks Release Pipeline' under 'Release-11'. The left sidebar includes options like Overview, Boards, Repos, Pipelines, Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The main area displays the 'Deployment process' with a status of 'Succeeded'. A specific 'Agent job' is highlighted, showing its tasks and their outcomes:

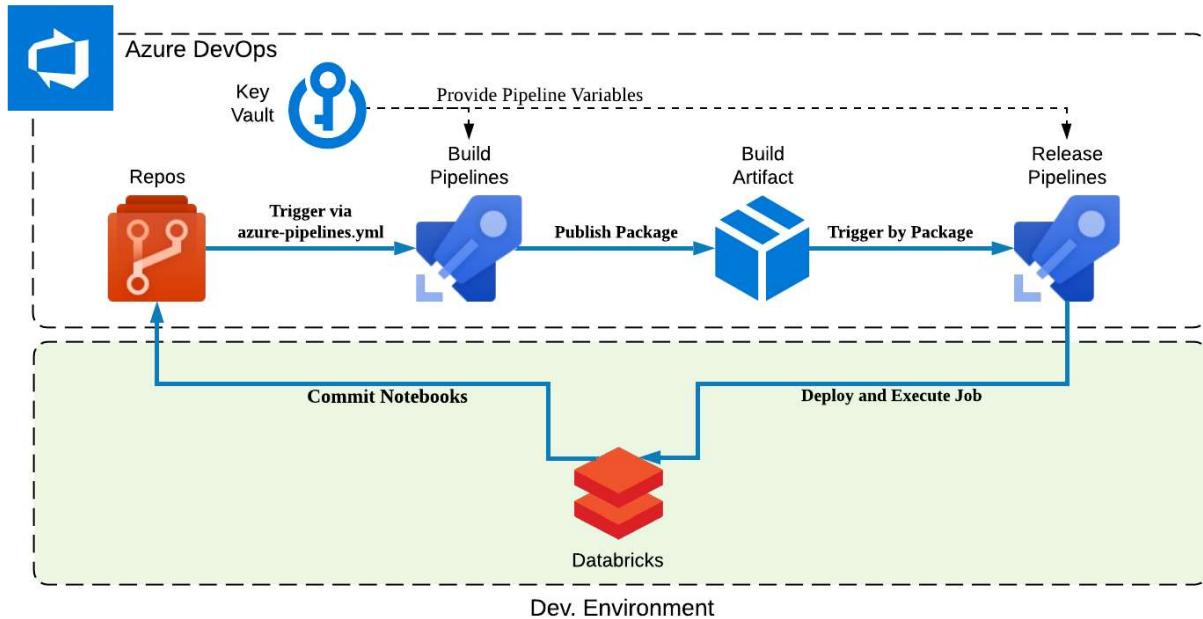
| Task | Status | Duration |
|--|-----------|----------|
| Initialize Agent | succeeded | <1s |
| Initialize job | succeeded | 4s |
| Download Artifacts | succeeded | <1s |
| Download secrets: bpdo-databricks-kv-dev | succeeded | 2s |
| Use Python 3.x | succeeded | <1s |
| Create Databricks Resource Group | succeeded | 3m 20s |
| Install Tools | succeeded | 12s |
| Authenticate with Databricks CLI | succeeded | <1s |
| Upload Notebook to Databricks | succeeded | 5s |
| Create Notebook Run JSON | succeeded | <1s |
| Run Notebook on Databricks | succeeded | 2s |
| Wait for Databricks Run to complete | succeeded | 4m 58s |

4–10 Deployment and execution result

Tada! Now your notebook is being deployed back to your Development environment, and successfully executed via a Job cluster!

Why do we want to do that? It is because you would like to test if the notebook can be executed on a cluster other than the interactive cluster you have been developing your notebook, ensuring your notebook is portable.

Let's recap what we have done



1. We have set up Databricks and Azure Key Vault provisioning via Azure Resource Manager Template.
2. We have set up Git integration with Databricks.
3. We have set up preliminary build steps and publish notebook as build artifacts.
4. We have been using Azure Key Vault for securely managing deployment credentials.
5. We have set up automated deployment and job execution flow with Databricks, which the job execution can be served as a very simple deployment test.

Edit on 3 March 2019

We have found that the key vault cannot be created automatically with the current ARM template, this tutorial will be updated in future. For now, please manually create the key vault for the first time in the resource group.

What's Next?

In our next posts, we will add Data Factory into the picture, provisioning sandbox data factory via ARM Template, and setting up Data Factory CI/CD properly.

Stay tuned!

Azure Databricks DevOps Ci Cd Pipeline Databricks Azure Key Vault

About Write Help Legal

Get the Medium app

