

DATABRICKS HANDS-ON LAB

Overview	3
Lab 1: Create the Data Factory and Load the Data to an Azure Blob Storage	4
Overview	4
Create An Azure Blob Storage.....	4
Create an Azure Data Factory and Load files from a remote Storage to our Blob Storage	7
Create your Databricks cluster	14
Terms of use	17

OVERVIEW

Most companies already have one or more data warehouses. However, extending and maintaining this data warehouse can be difficult. Source systems are changing faster than ever before, and end users want to make deeper analyses.

Therefore, a more flexible architecture is needed which makes it easier to add different types of data.

During this workshop you will experience extend the data warehouse using the Azure Data Services.

The use case during this workshop is about airdelays and preparing the data for Data Scientists on the one hand but also providing it for analysts via the Data Warehouse.

This Lab will guide through the data acquisition and how to create data pipelines with Azure Data Factory and load data into an Azure Blob Storage for further usage with Azure Databricks.

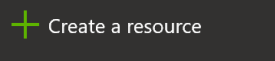
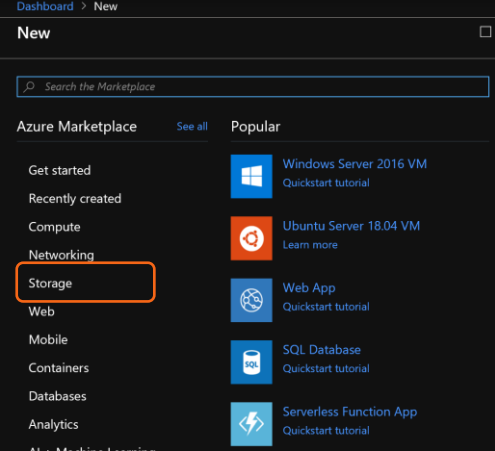
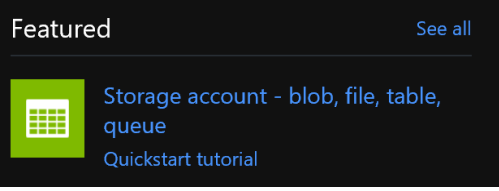
Good luck and enjoy the work!

LAB 1: CREATE THE DATA FACTORY AND LOAD THE DATA TO AN AZURE BLOB STORAGE

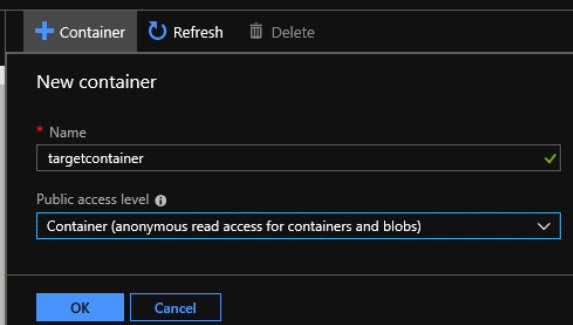
OVERVIEW

This first lab for today will walk you through the creation of a Azure Blob Storage and the use of Azure Data Factory to fetch files from a folder on Azure, that we then will use as the data in all the subsequent labs.

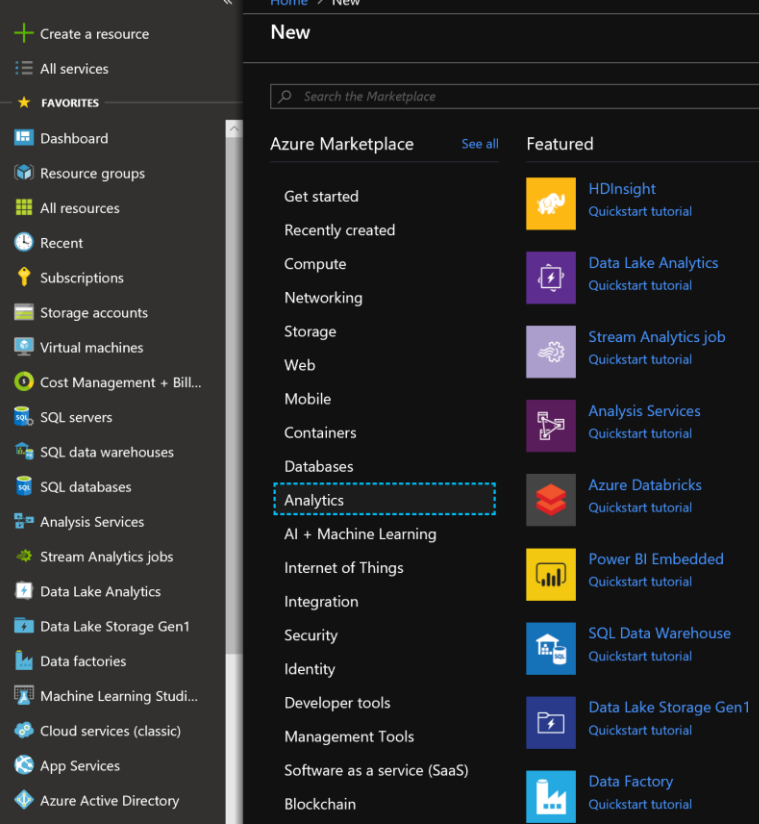
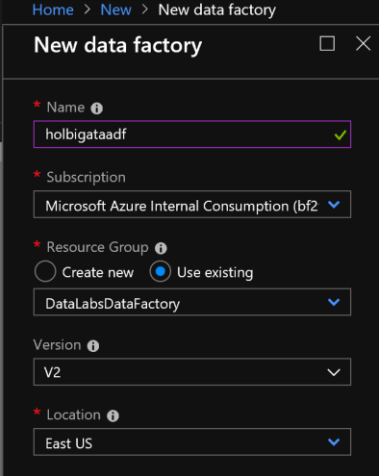
CREATE AN AZURE BLOB STORAGE

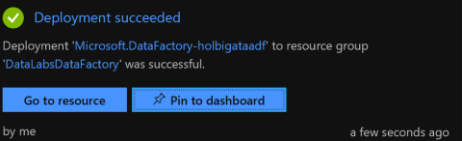
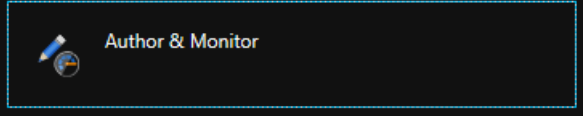
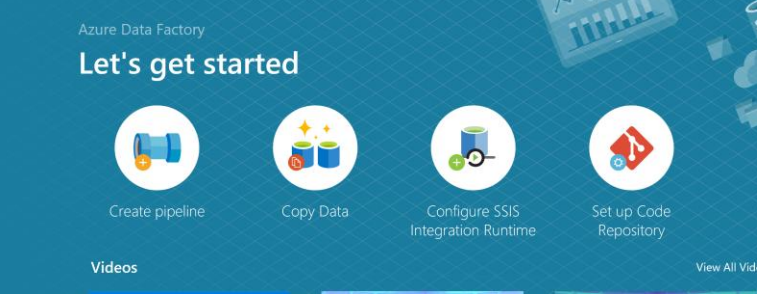
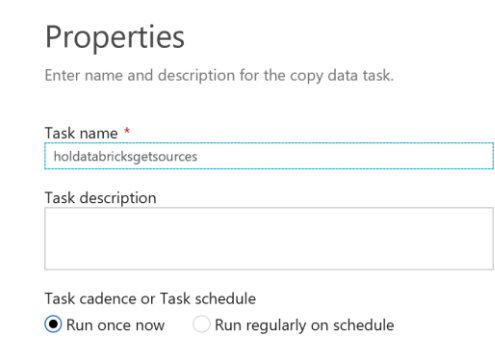
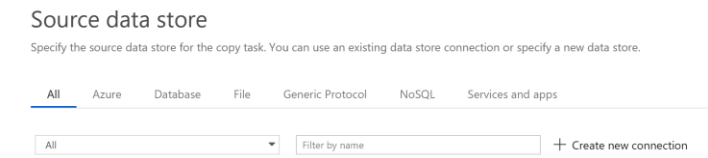
	On the Azure Portal please search for the 'Create a resource' button.
	Within the opening blade you will find the entry for 'Storage'
	Please go for 'Storage account – blob, ...'

	<p>In the first drop-down box please search for your subscription.</p> <p>In the second box please create a new resource group. This group can be used throughout the whole training. Afterwards you can delete all created artifacts by just deleting this resource group.</p> <p>Please provide a 'Storage account name'. This name must be unique over all the storage accounts in Azure and can only contain lower letters and numbers.</p> <p>When choosing the location for the service please take care, that all services created today are in the same region to avoid additional data egress cost.</p> <p>Leave the rest of the properties as is and hit 'Next: Advanced >'</p>
	<p>On this blade, leave all entries on their default values and hit 'Review + create' → the storage account will be created for you.</p>
	<p>After the storage account is created, please open it and select 'Blobs'</p>
	<p>Please click on '+ Container' to create a new container, to store the incoming data.</p>

	<p>Name it and then adjust the 'Public access level' like in the screenshot and hit 'OK'.</p> <p>This will be the target folder for the next step.</p>
---	--

CREATE AN AZURE DATA FACTORY AND LOAD FILES FROM A REMOTE STORAGE TO OUR BLOB STORAGE


	<p>Go to your Azure portal and click +, choose Analytics then choose Data Factory</p>
	<p>Enter a unique name.</p> <p>Choose the resource group you have created in the last sequence for your storage service</p> <p>You want to use Version 'V2' for this lab.</p> <p>And don't forget to put the Data Factory into the region, that you chose for the storage account above. (it will definitely differ from the one in the picture 😊)</p>

	<p>Hit 'create' and wait for the service to be displayed. After creation, the Bell Sign will show new messages.</p> <p>Click on it, choose the Deployment-Success message for your Data Factory and click 'Pin to dashboard'.</p>
	<p>Now open your Data Factory: search for 'Author and Monitor'</p>
	<p>and in the new tab (takes a short amount of time) start the editor: 'Copy Data'.</p>
	<p>On the first screen name your data copy pipeline like 'holdatabricksgetsourcedata' or just leave the default name and click 'Next'</p>
	<p>Here we will create a new data store, that will be used as a connection within Data Factory Pipelines. In this case it'll be the source data for your pipeline.</p>


New Linked Service

Search


All Azure Database File Generic Protocol NoSQL Services and apps




Amazon Marketplace Web Service (Preview)




Amazon Redshift




Amazon S3




Apache Impala (Preview)




Azure Blob Storage




Azure Cosmos DB



Azure Data Lake Storage Gen1



Azure Data Lake Storage Gen2 (Preview)



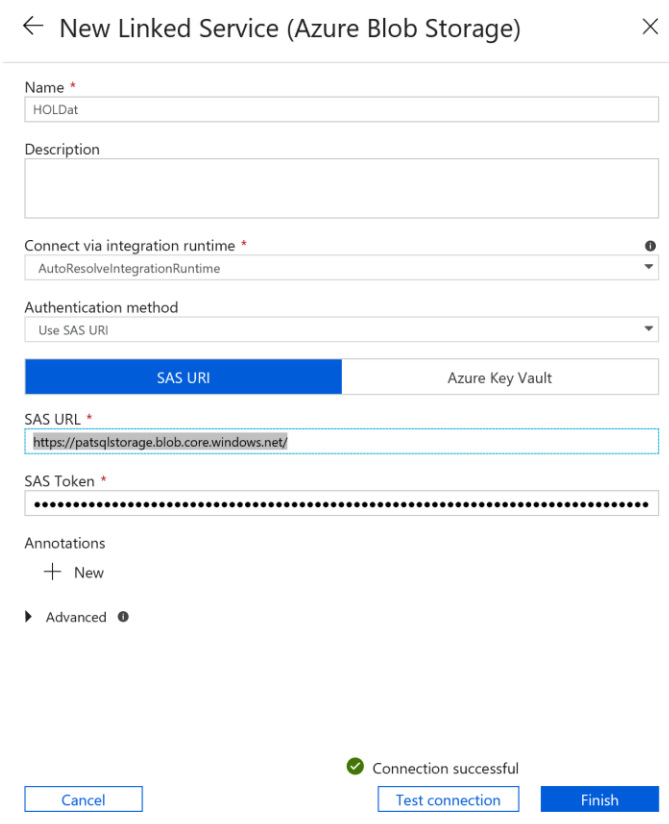
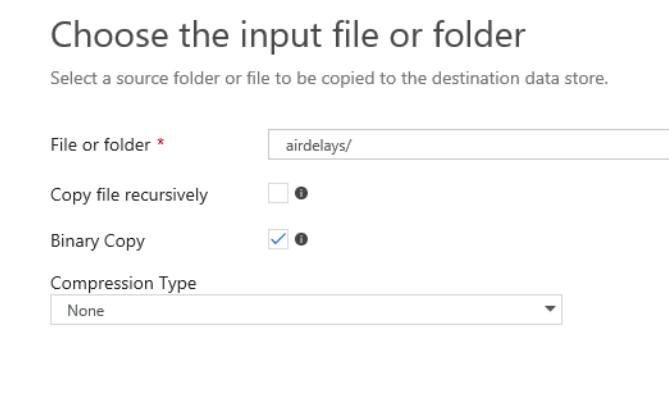
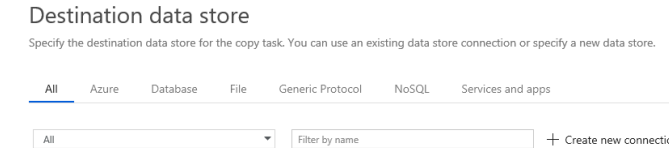
Azure Database for MariaDB

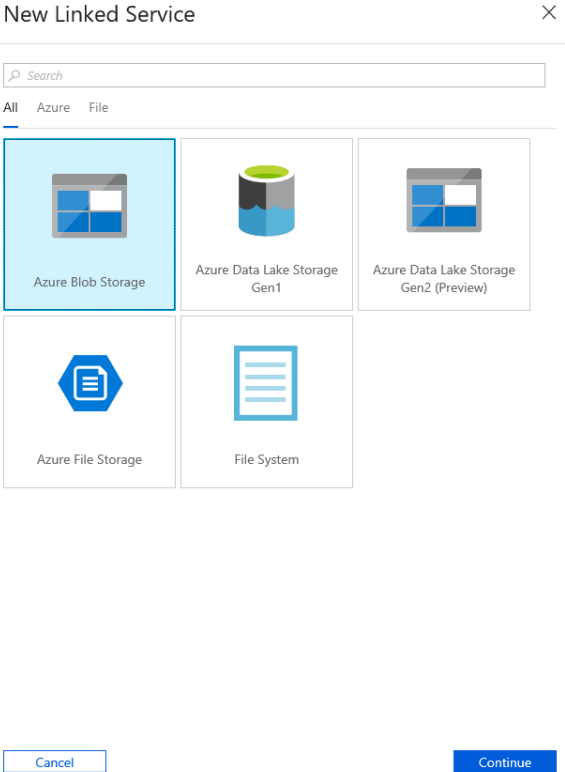
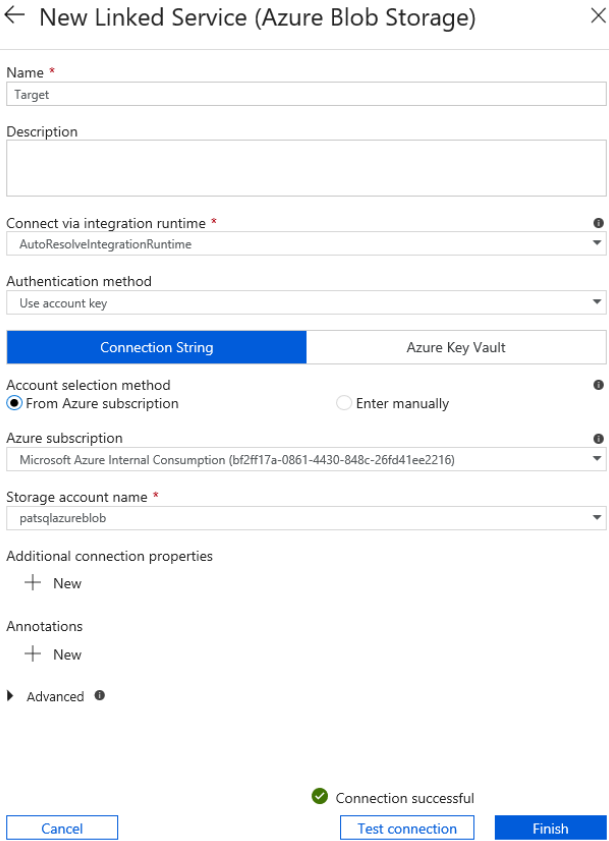
Cancel

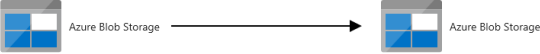
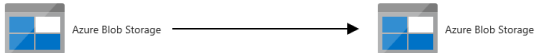
Continue

Click on 'Azure Blob Storage' and hit 'Continue'

9

 <p>← New Linked Service (Azure Blob Storage) ×</p> <p>Name * HOLDat</p> <p>Description</p> <p>Connect via integration runtime * AutoResolveIntegrationRuntime</p> <p>Authentication method Use SAS URI</p> <p>SAS URI Azure Key Vault</p> <p>SAS URL * https://patsqlstorage.blob.core.windows.net/</p> <p>SAS Token *</p> <p>Annotations + New</p> <p>▶ Advanced ⓘ</p> <p>Connection successful</p> <p>Cancel Test connection Finish</p>	<p>On this dialog please name your connection.</p> <p>Please adjust the 'Authentication method' to 'Use SAS URI'.</p> <p>For the SAS URL please use: https://patsqlstorage.blob.core.windows.net/</p> <p>and for the SAS Token please use: VCXxoTmA03iMVAY+IOHtVoYXvpsB8W Fni1/QcqEYjCAAUComxLDTyiQ65t3Ag6 1FBoypsbqgAcRVizsSUaFGog==</p> <p>You don't need to add other information and can click on 'Test connection' to check connectivity. If the connection turns green, you may hit 'Finish'.</p> <p>With the newly created source marked, please hit 'Next'.</p>
 <p>Choose the input file or folder</p> <p>Select a source folder or file to be copied to the destination data store.</p> <p>File or folder * airdelays/</p> <p>Copy file recursively <input type="checkbox"/> ⓘ</p> <p>Binary Copy <input checked="" type="checkbox"/> ⓘ</p> <p>Compression Type None</p>	<p>On the following screen please select the source folder → click 'Browse', select the folder and click 'Choose'.</p> <p>Then please check the 'Binary Copy' property.</p> <p>Hit 'Next'</p>
 <p>Destination data store</p> <p>Specify the destination data store for the copy task. You can use an existing data store connection or specify a new data store.</p> <p>All Azure Database File Generic Protocol NoSQL Services and apps</p> <p>All Filter by name + Create new connection</p>	<p>On the 'Destination data store' screen please again create another connection, this time to the storage account you have created above.</p>

	<p>Please create a new Blob Storage connection.</p>
	<p>Name the target connection.</p> <p>Please leave 'Use account key' in the 'Authentication method' box.</p> <p>As the 'Account selection method' please also leave 'From Azure subscription'.</p> <p>In the 'Azure subscription' selector please pick your subscription.</p> <p>In the 'Storage account name' box select the account, you have created above.</p> <p>If the connection tests successful, please hit 'Finish' and with the newly created source selected click 'Continue' and proceed to the folder picker dialog.</p>

<h3>Choose the output file or folder</h3> <p>Specify a folder that will contain output files or a specific output file in the destination data store.</p> <p>Folder path * <input type="text" value="tmpairdelay"/></p> <p>File name <input type="text" value="File names are defined by source"/></p> <p>Compression Type <input type="text" value="None"/></p> <p>Copy behavior <input type="text" value="Preserve hierarchy"/></p>	<p>Choose the folder created above and hit 'Next'</p>
<h3>Settings</h3> <p>More options for data movement</p> <p>▸ Performance settings</p> <p>Enable Staging <input type="checkbox"/></p> <p>▸ Advanced settings</p>	<p>This dialog can be left like it is. (if you would like to stage data, this would be the place to configure a folder for staging files before they're loaded to another Azure Data Service).</p>
<h3>Summary</h3> <p>You are running pipeline to copy data from Azure Blob Storage to Azure Blob Storage.</p>  <p>Properties Edit</p> <p>Task name CopyPipeline_g1v</p> <p>Task description</p> <p>Source Edit</p> <p>Connection name AzureBlobStorage1</p> <p>Dataset name SourceDataset_g1v</p> <p>File name</p> <p>Directory path airdelays</p> <p>Destination Edit</p> <p>Connection name Target</p> <p>Dataset name DestinationDataset_g1v</p> <p>File name</p> <p>Directory path tmpairdelay</p> <p>Copy settings Edit</p> <p>Timeout 7 00:00:00</p> <p>Previous Next</p>	<p>On the 'Summary' blade you will have the last chance to check and change settings.</p>
 <p>Deployment complete</p> <ul style="list-style-type: none"> ▸ Creating Datasets ✓ ▸ Creating Pipelines ✓ ▸ Running Pipelines ✓ <p>Datasets and pipelines have been created. You can now monitor and edit the copy pipelines or click finish to close the copy wizard.</p> <p>Edit Pipeline Monitor</p>	<p>After hitting 'Next' the pipeline will be deployed to the Data Factory. You will have the chance to monitor the pipeline execution.</p>

The screenshot shows the Databricks tmpairdelay interface. On the left is a sidebar with navigation options: Overview, Access Control (policies), Settings, Access policy, Properties, Metadata, and Editor (preview). The main panel displays a table titled 'Locations tmpairdelay' with a search bar and a 'Show deleted blobs' toggle. The table contains the following data:

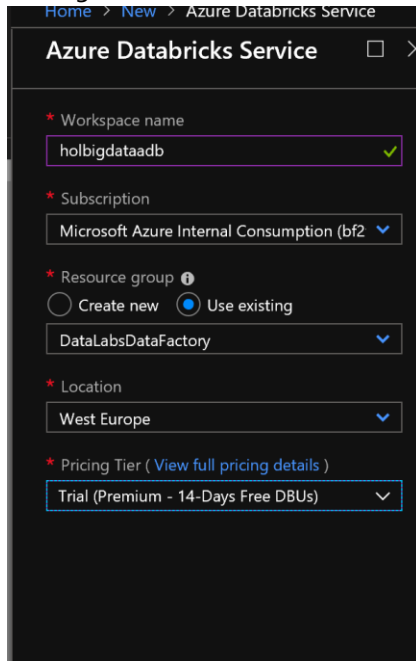
name	timestamp	access time	blob type	size	status
zaoT78076.csv	13/12/2018 5:30:11 PM		Block Blob	81.49 KB	Available
zaoT78076.csv	13/12/2018 5:29:57 PM		Block Blob	76.75 KB	Available
zaoT78076.csv	13/12/2018 5:30:03 PM		Block Blob	79.9 KB	Available
zaoT78081.csv	13/12/2018 5:30:10 PM		Block Blob	78.56 KB	Available
zaoT78082.csv	13/12/2018 5:29:52 PM		Block Blob	74.45 KB	Available
zaoT78083.csv	13/12/2018 5:30:11 PM		Block Blob	80.49 KB	Available
zaoT78084.csv	13/12/2018 5:30:08 PM		Block Blob	72.33 KB	Available

Now you can proceed to the target folder and check the data, that was loaded. It is now available for the processing with Databricks.

CREATE YOUR DATABRICKS CLUSTER

In the next step you will create an Azure Databricks cluster for processing the data that you have loaded.

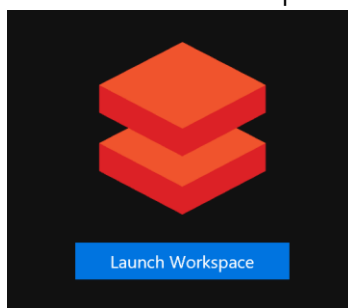
1. Please go to the portal, click '+ Create a resource' and on the Analytics-Tab search for 'Azure Databricks'
2. On the following dialogue please enter a unique name for your Databricks Cluster, the resource group that you have created in the first lab and select the region 'Western Europe' and the Pricing Tier:



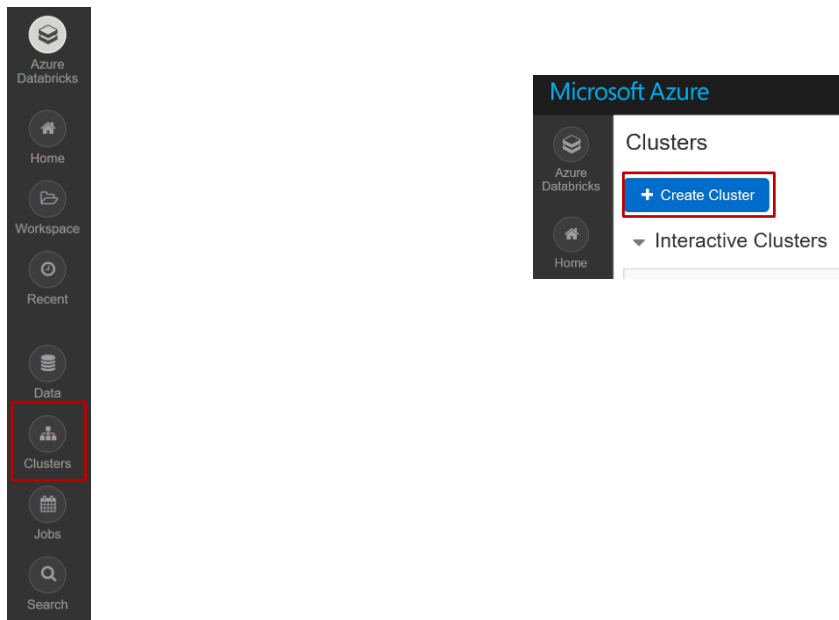
The screenshot shows the 'Azure Databricks Service' configuration window. The breadcrumb path at the top is 'Home > New > Azure Databricks Service'. The window title is 'Azure Databricks Service'. The configuration fields are as follows:

- Workspace name:** 'holbigdataadb' with a green checkmark.
- Subscription:** 'Microsoft Azure Internal Consumption (bf2)' with a dropdown arrow.
- Resource group:** 'DataLabsDataFactory' with a dropdown arrow. The options are 'Create new' (radio button) and 'Use existing' (radio button, selected).
- Location:** 'West Europe' with a dropdown arrow.
- Pricing Tier:** 'Trial (Premium - 14-Days Free DBUs)' with a dropdown arrow. A link '(View full pricing details)' is visible.

3. After the Databricks workspace is created and shows up on the dashboard please select it and click on 'Launch Workspace'. It will show up in a different tab:



4. As the workspace launches we can start and create the first Databricks cluster. Please click 'Clusters' and then on the next screen 'Create Cluster'



5. On the following screen you will be able to configure the planned cluster as you need it. Please enter a Cluster Name, select the Databricks Runtime Version (leave the default), select the Python Version = 3 and the size of the Cluster-VMs (Driver and Worker) and hit 'Create Cluster'

Create Cluster

New Cluster Cancel Create Cluster 2-8 Workers: 28.0-112.0 GB Memory, 8-32 Cores, 1.5-6 DBU
1 Driver: 14.0 GB Memory, 4 Cores, 0.75 DBU Cost \$0.40 per DBU

Cluster Name Please enter a cluster name

Cluster Mode
☐ High Concurrency Optimized to run concurrent SQL, Python, and R workloads. Does not support Scala. Previously known as Serverless.
☒ Standard Recommended for single-user clusters. Can run SQL, Python, R, and Scala workloads.

Databricks Runtime Version ?
4.2 (includes Apache Spark 2.3.1, Scala 2.11)

Python Version ?
2

Driver Type
Same as worker 14.0 GB Memory, 4 Cores, 0.75 DBU

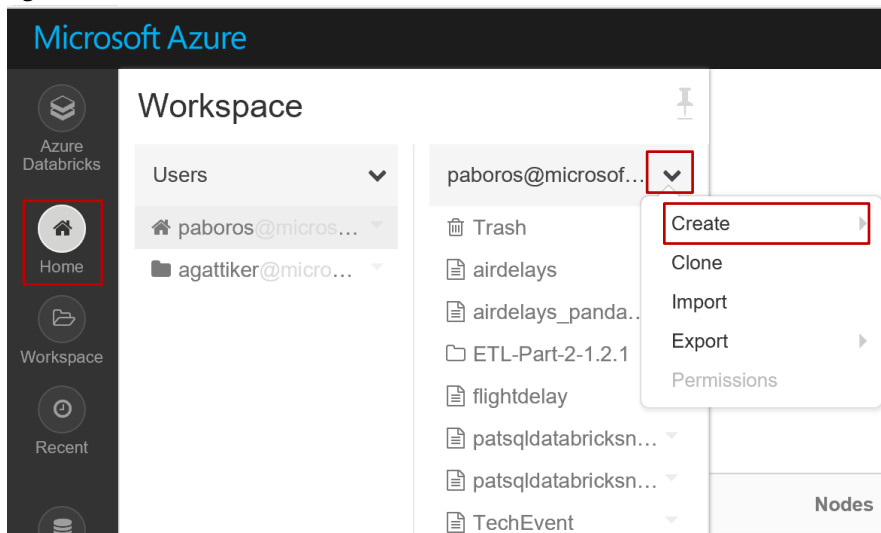
Worker Type 14.0 GB Memory, 4 Cores, 0.75 DBU
Standard_DS3_v2

Min Workers 2 Max Workers 8 ☒ Enable autoscaling ?

Auto Termination ?
☒ Terminate after 120 minutes of inactivity

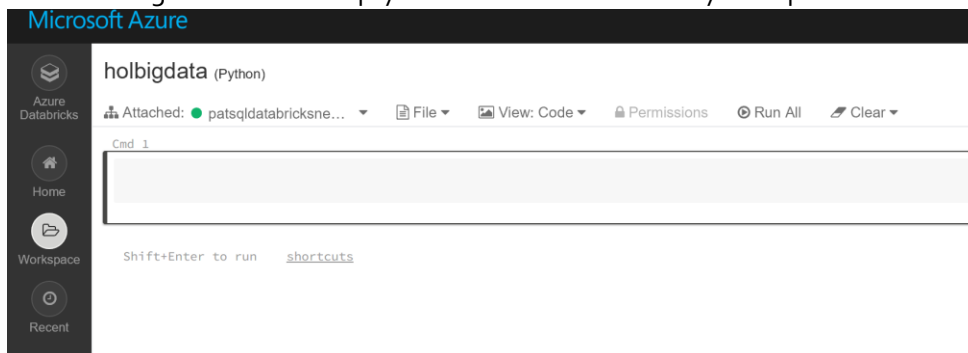
Spark [Tags](#) [Logging](#) [Init Scripts](#)

6. To use the Cluster we need to create a notebook now, where we can enter code and run it against it. Please click on 'Home' and on the Users tab on the arrow. Then select 'Create'



7. On the next cartridge select 'Notebook'. Name the new notebook. The rest can be left defaulted

8. After clicking 'Create' the empty notebook is available for your input.



TERMS OF USE

© 2017 Microsoft Corporation. All rights reserved.

By using this Hands-on Lab, you agree to the following terms:

The technology/functionality described in this Hands-on Lab is provided by Microsoft Corporation in a “sandbox” testing environment for purposes of obtaining your feedback and to provide you with a learning experience. You may only use the Hands-on Lab to evaluate such technology features and functionality and provide feedback to Microsoft. You may not use it for any other purpose. You may not modify, copy, distribute, transmit, display, perform, reproduce, publish, license, create derivative works from, transfer, or sell this Hands-on Lab or any portion thereof.

COPYING OR REPRODUCTION OF THE HANDS-ON LAB (OR ANY PORTION OF IT) TO ANY OTHER SERVER OR LOCATION FOR FURTHER REPRODUCTION OR REDISTRIBUTION IS EXPRESSLY PROHIBITED.

THIS HANDS-ON LAB PROVIDES CERTAIN SOFTWARE TECHNOLOGY/PRODUCT FEATURES AND FUNCTIONALITY, INCLUDING POTENTIAL NEW FEATURES AND CONCEPTS, IN A SIMULATED ENVIRONMENT WITHOUT COMPLEX SET-UP OR INSTALLATION FOR THE PURPOSE DESCRIBED ABOVE. THE TECHNOLOGY/CONCEPTS REPRESENTED IN THIS HANDS-ON LAB MAY NOT REPRESENT FULL FEATURE FUNCTIONALITY AND MAY NOT WORK THE WAY A FINAL VERSION MAY WORK. WE ALSO MAY NOT RELEASE A FINAL VERSION OF SUCH FEATURES OR CONCEPTS. YOUR EXPERIENCE WITH USING SUCH FEATURES AND FUNCTIONALITY IN A PHYSICAL ENVIRONMENT MAY ALSO BE DIFFERENT.

FEEDBACK. If you give feedback about the technology features, functionality and/or concepts described in this Hands-on Lab to Microsoft, you give to Microsoft, without charge, the right to use, share and commercialize your feedback in any way and for any purpose. You also give to third parties, without charge, any patent rights needed for their products, technologies and services to use or interface with any specific parts of a Microsoft software or service that includes the feedback. You will not give feedback that is subject to a license that requires Microsoft to license its software or documentation to third parties because we include your feedback in them. These rights survive this agreement.

MICROSOFT CORPORATION HEREBY DISCLAIMS ALL WARRANTIES AND CONDITIONS WITH REGARD TO THE HANDS-ON LAB, INCLUDING ALL WARRANTIES AND CONDITIONS OF MERCHANTABILITY, WHETHER EXPRESS, IMPLIED OR STATUTORY, FITNESS FOR A PARTICULAR PURPOSE, TITLE AND NON-INFRINGEMENT. MICROSOFT DOES NOT MAKE ANY ASSURANCES OR REPRESENTATIONS WITH REGARD TO THE ACCURACY OF THE RESULTS, OUTPUT THAT DERIVES FROM USE OF THE VIRTUAL LAB, OR SUITABILITY OF THE INFORMATION CONTAINED IN THE VIRTUAL LAB FOR ANY PURPOSE.

DISCLAIMER

This lab contains only a portion of the features and enhancements in Microsoft Azure Data Factory, Azure SQL Data Warehouse, Azure DataBrics and Azure Data Lake Storage. Some of the features might change in future releases of the product.