# Deploying HDInsight Spark Cluster on Azure

Step by Step

Basics  Storage  Security + networking  Configuration + pricing  Tags  Review + create

New to HDInsight? Get started with our training resources.
Create a managed HDInsight cluster. Select from Spark, Kafka, Hadoop, Storm, and more. Learn more

**Project details**

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *           springboard-data-engineering

Resource group *        SparkLabResourceGroup
Create new

**Cluster details**

Name your cluster, pick a region, and choose a cluster type and version. Learn more

Cluster name *           sparklabjwittbold

Region *                 Southeast Asia

Cluster type *           **Spark**
                         Change

Version *                Spark 2.4 (HDI 4.0)

**Cluster credentials**

Enter new credentials that will be used to administer or access the cluster.

Cluster login username *  ⓘ    admin

Cluster login password *        ••••••••••••

Confirm cluster login password *  ••••••••••••

Secure Shell (SSH) username *  ⓘ   sshuser

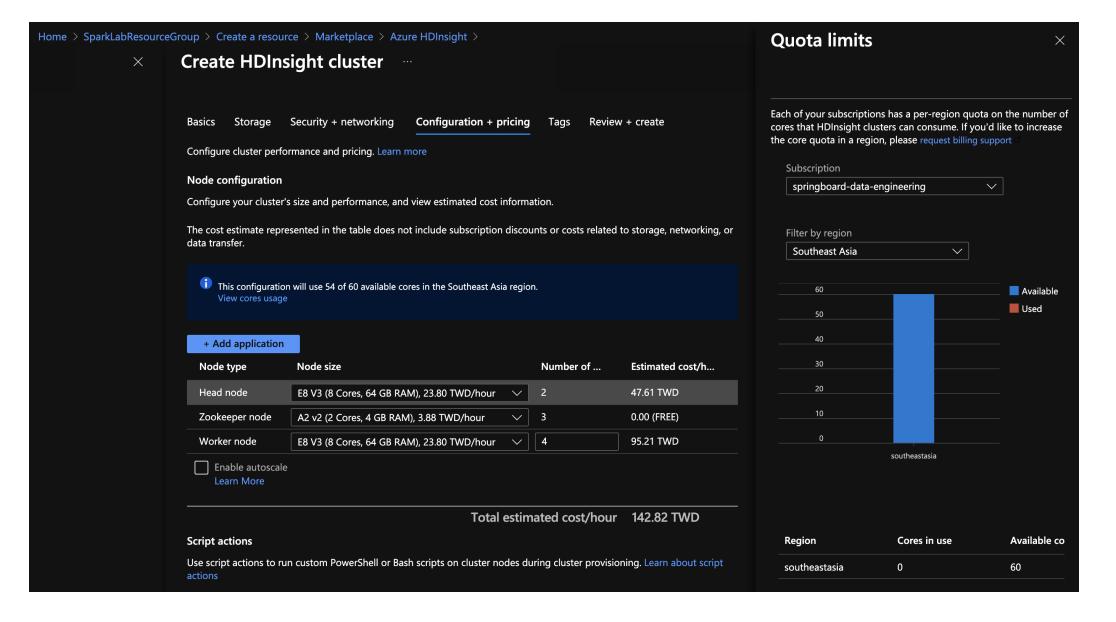Use cluster login password for SSH   ☑

- Select Subscription
- Select Resource Group

- Provide Cluster name
- Choose Region
- Specify Cluster Type / Version

- Set password to login as admin over SSH

Select or create storage accounts that will be used for the cluster's logs, job input, and job output. Configure the cluster's access to these accounts, if needed.

**Primary storage**

Select or create a storage account that will be the default location for cluster logs and other output.

Primary storage type *                    | Azure Storage                                          ▽ |

Selection method * ⓘ          ◉ Select from list      ◯ Use access key

Primary storage account *          | sparklabstoragejwittbold                               ▽ |

Create new

Container * ⓘ                       | sparklabjwittbold-container                            ✓ |

**Data Lake Storage Gen1**

Provide details for the cluster to access Data Lake Storage Gen1. The cluster will be able to access any Data Lake Storage Gen1 accounts that the chosen service principal has access to.

Data Lake Storage Gen1 access          Configure access settings

**Additional Azure Storage**

Link additional Azure Storage accounts to the cluster.

Add Azure Storage

**Custom Ambari DB**

Use an external Ambari database for greater flexibility, control, and customization. Learn More

SQL database for Ambari ⓘ          |                                                        ▽ |

**External metadata stores**

To store your Hive and Oozie metadata outside of this cluster, select a SQL database. Learn More

SQL database for Hive ⓘ          |                                                        ▽ |

SQL database for Oozie ⓘ          |                                                        ▽ |

- Choose Storage Account
- Set name for Container

Configure your cluster's security and network settings.

**Enterprise security package**

Connect this cluster with Active Directory Domain Services (AAD-DS) to have finer control of who can access the cluster. Learn More

☐ Enable enterprise security package ⓘ (Adds 0.3155691 TWD per Core-Hour)

**TLS**

Select the minimum TLS version supported for your cluster. Learn more

Minimum TLS version ⓘ

| 1.2 | ⌄ |
|---|---|

**Network settings**

Resource provider connection ⓘ

| Inbound | ⌄ |
|---|---|

Connect this cluster to a virtual network. Learn more

Virtual network ⓘ

| | ⌄ |
|---|---|

**Encryption in transit**

Configure encryption in transit settings. Learn more

☐ Enable encryption in transit ⓘ

**Encryption at rest**

Configure disk encryption settings. Learn more

☐ Provide your own key from key vault ⓘ

☐ Enable encryption at host on temp data disk ⓘ

**Identity**

Select a user-assigned service identity to represent your cluster for enterprise security package or disk encryption. Learn more

User-assigned managed identity ⓘ

| | ⌄ |
|---|---|

Configure Security
(not modified)

Configure Nodes
- Needed to first request Core Quota increase for subscription / region

# Create HDInsight cluster ···

✅ Validation succeeded.

Basics    Storage    Security + networking    Configuration + pricing    Tags    **Review + create**

Spark 2.4 (HDI 4.0)        **142.82 TWD Total estimated cost/hour**
                          This estimate does not include subscription discounts or costs related to storage,
                          networking, or data transfer.

## Basics

| | |
|---|---|
| Subscription | springboard-data-engineering |
| Resource group | SparkLabResourceGroup |
| Region | Southeast Asia |
| Cluster name | (new) sparklabjwittbold |
| Cluster type | Spark 2.4 (HDI 4.0) |
| Cluster login username | admin |
| Secure Shell (SSH) username | sshuser |
| Use cluster login password for SSH | Enabled |

## Security + networking

| | |
|---|---|
| Minimum TLS version | 1.2 |
| Resource provider connection | Inbound |
| Encryption at rest | Disabled |
| Encryption in transit | Disabled |
| Encryption at host on temp data disk | Disabled |

## Storage

| | |
|---|---|
| Primary storage type | Azure Storage |
| Primary storage account | sparklabstoragejwittbold |
| Container | sparklabjwittbold-container |
| Additional Azure Storage | None |
| Data Lake Storage Gen1 access | Disabled |

## Cluster configuration
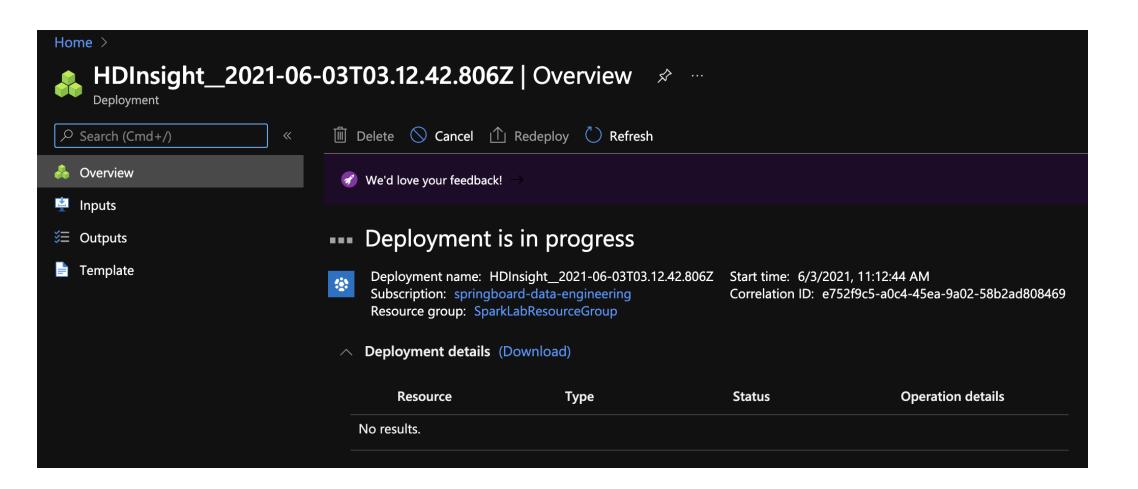
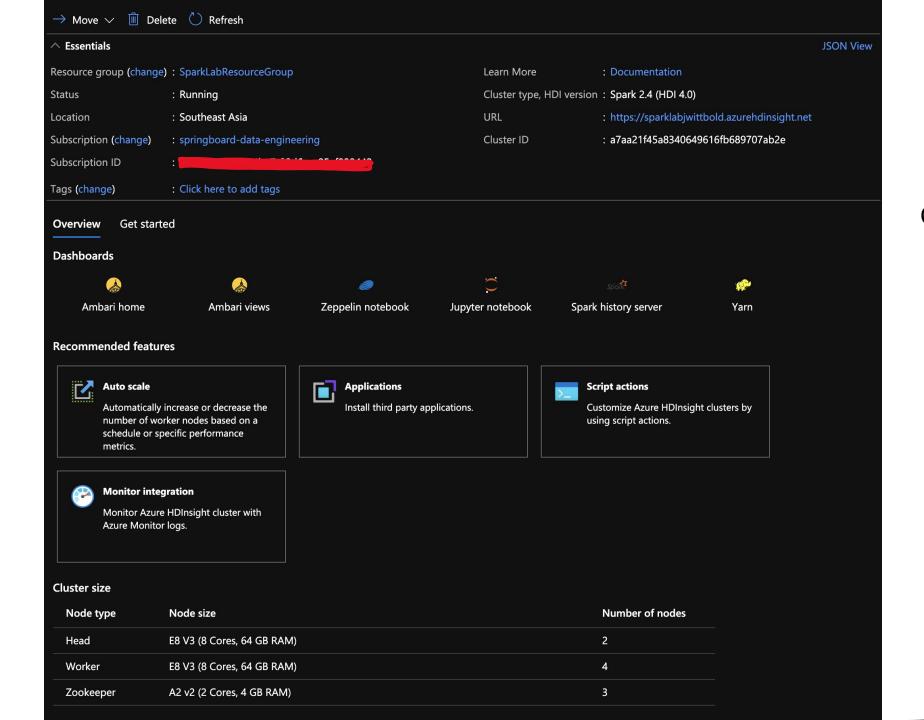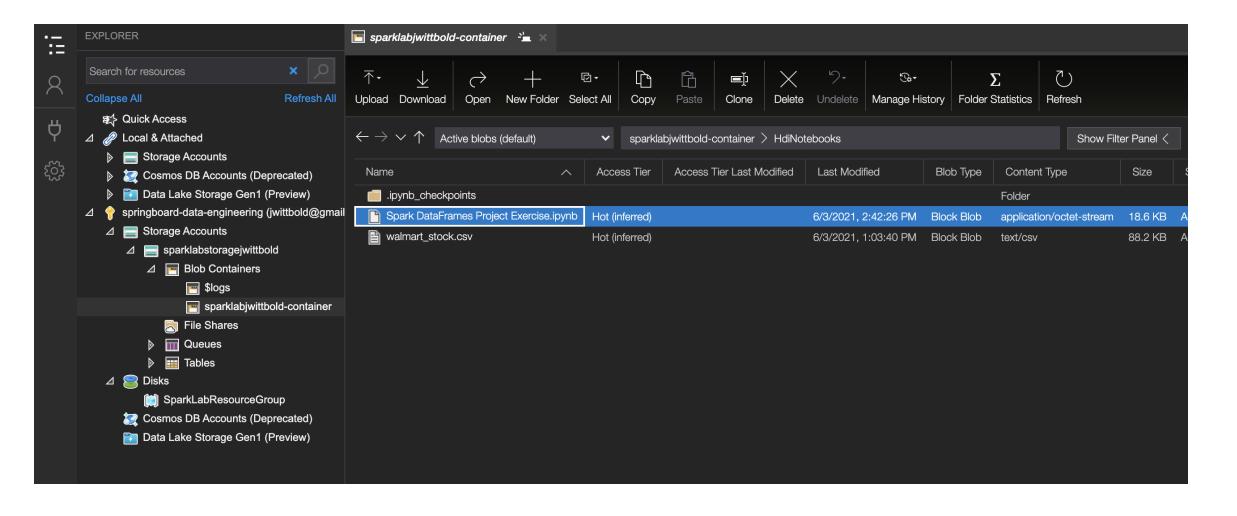| | |
|---|---|
| Head | 2 nodes, E8 V3 (8 Cores, 64 GB RAM) |
| Zookeeper | 3 nodes, A2 v2 (2 Cores, 4 GB RAM) |
| Worker | 4 nodes, E8 V3 (8 Cores, 64 GB RAM) |

Full Details of HDInsight Cluster
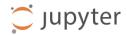
- Validation Successful

Deployment in progress...

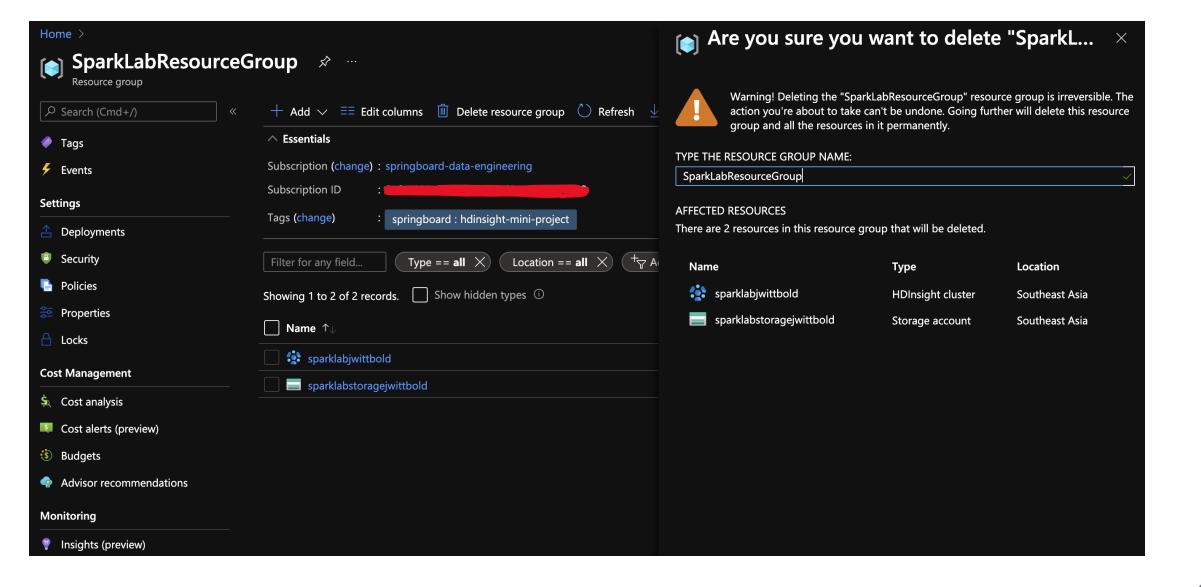⌃ **Essentials**                                                                    JSON View

Resource group (change)    : SparkLabResourceGroup          Learn More            : Documentation

Status                     : Running                        Cluster type, HDI version : Spark 2.4 (HDI 4.0)

Location                   : Southeast Asia                 URL                   : https://sparklabjwittbold.azurehdinsight.net

Subscription (change)      : springboard-data-engineering   Cluster ID            : a7aa21f45a8340649616fb689707ab2e

Subscription ID            : ▬▬▬▬▬▬▬▬▬▬▬

Tags (change)              : Click here to add tags

**Overview**   Get started

**Dashboards**

🟠                  🟠                  🔵                  ⚪                  spark⭐️             🐘
Ambari home        Ambari views       Zeppelin notebook  Jupyter notebook   Spark history server  Yarn

**Recommended features**

| ⬀ **Auto scale** | ⬛ **Applications** | ⌗ **Script actions** |
|---|---|---|
| Automatically increase or decrease the number of worker nodes based on a schedule or specific performance metrics. | Install third party applications. | Customize Azure HDInsight clusters by using script actions. |

| ⏲ **Monitor integration** | | |
|---|---|---|
| Monitor Azure HDInsight cluster with Azure Monitor logs. | | |

**Cluster size**

| Node type | Node size | Number of nodes |
|---|---|---|
| Head | E8 V3 (8 Cores, 64 GB RAM) | 2 |
| Worker | E8 V3 (8 Cores, 64 GB RAM) | 4 |
| Zookeeper | A2 v2 (2 Cores, 4 GB RAM) | 3 |

# Cluster successfully created

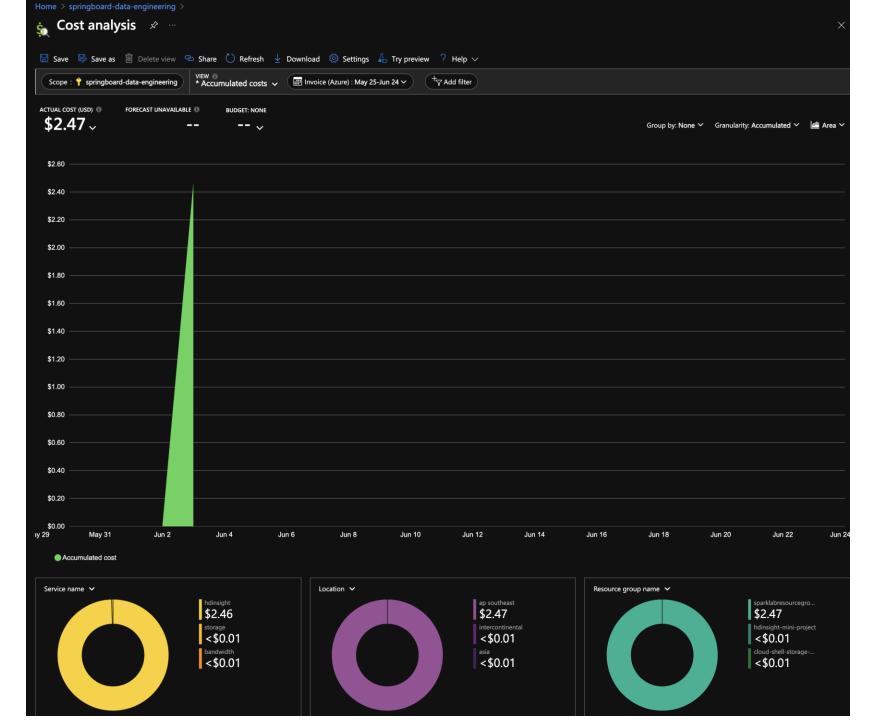Resource files uploaded to cluster via
Azure Storage Explorer

Jupyter Notebook and walmart_stock.csv data file available to work with

Deleting Resource Group (and all contained contents) after completing exercise

Successfully deleted

Cost Analysis for ~ 3hr
HDInsight Cluster

= $2.47 USD