# Azure Databricks Architecture and Security

databricks

# Quick Agenda

- Azure Databricks Platform Architecture
- Azure Databricks Security
  - Data Protection
  - IAM/Auth
  - Network Security
  - Compliance

# Azure Databricks Platform Architecture

# Azure Databricks Platform Architecture

## Microsoft Subscription

### Control Plane

| | |
|---|---|
| Web Application | Cluster Manager |
| Jobs | Notebooks |
| Hive Metastore | ACLs/Sessions |

### Backend Services
Log Storage/Analysis, Central Account Directory, Monitoring

## Client Subscription

### Data Plane

### Customer Data Sources

databricks

# Azure Databricks Platform Architecture Cont'd

## Microsoft/Databricks Subscription

## Client Subscription

**Control Plane**

**Data Plane**

Commands, table metadata, logs

Data being read from and written to the data source

User and organization meta-data, logs, useage

**Backend Services**

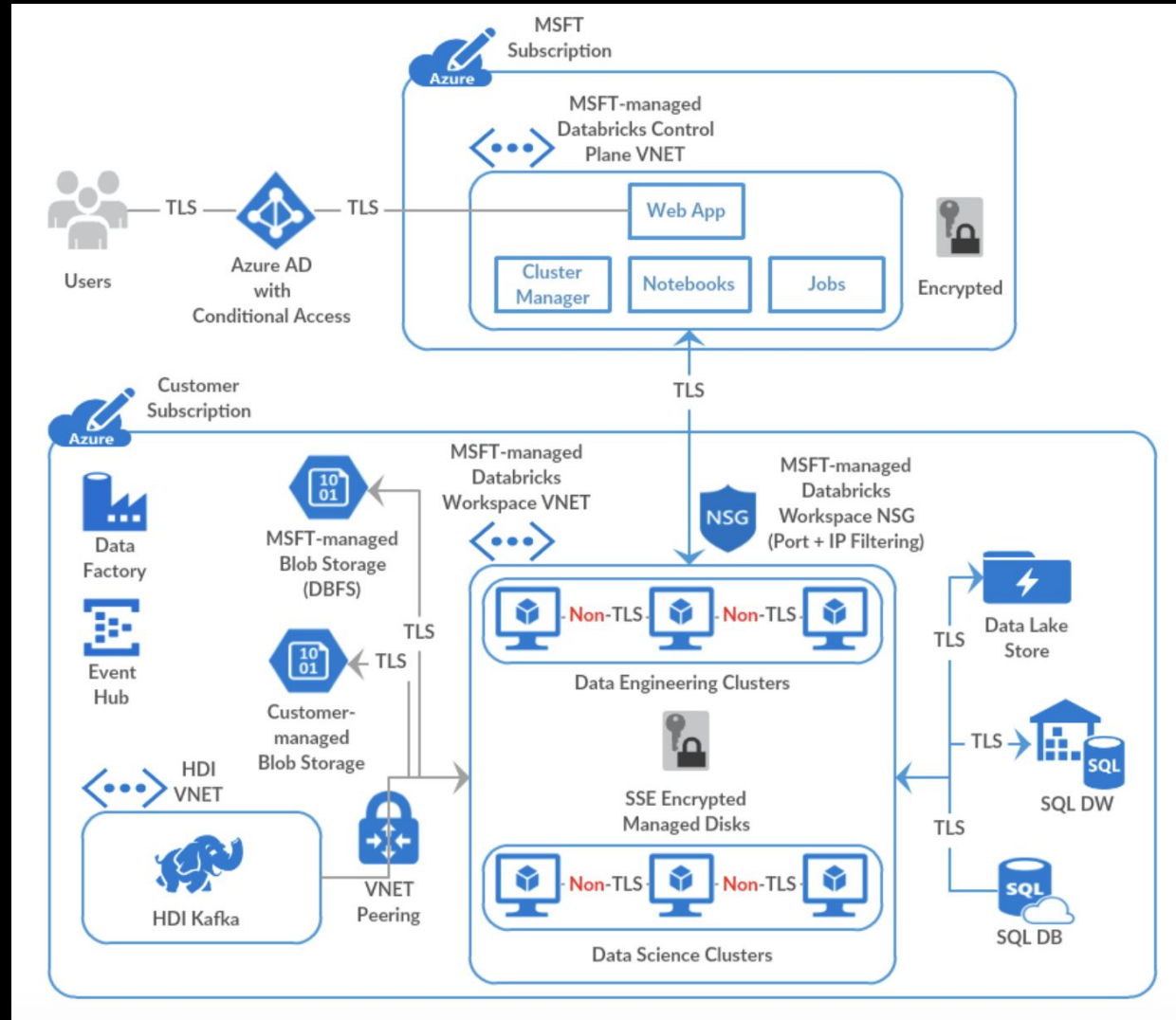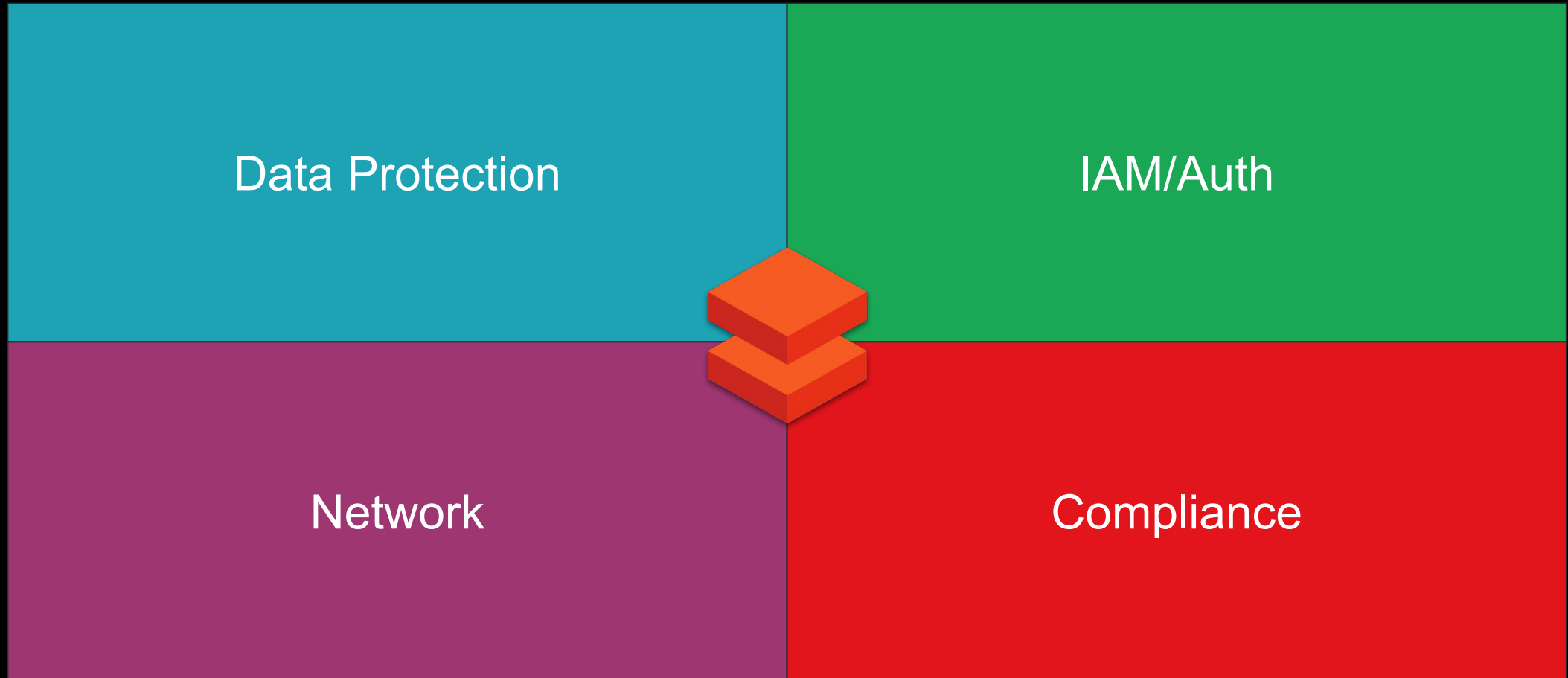**Customer Data Sources**

databricks

# Azure Databricks Platform Architecture

Standard Deployment View with no inter-node TLS

# Azure Databricks Security

# Azure Databricks Security

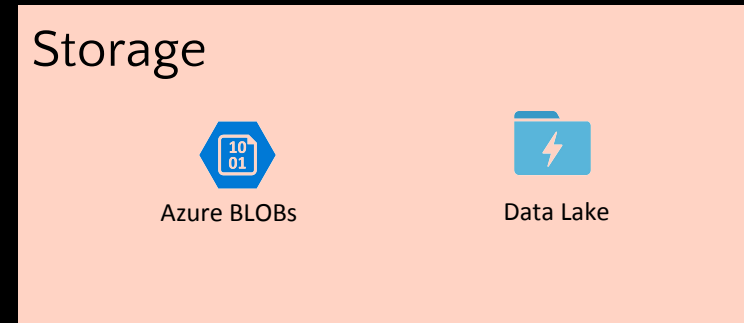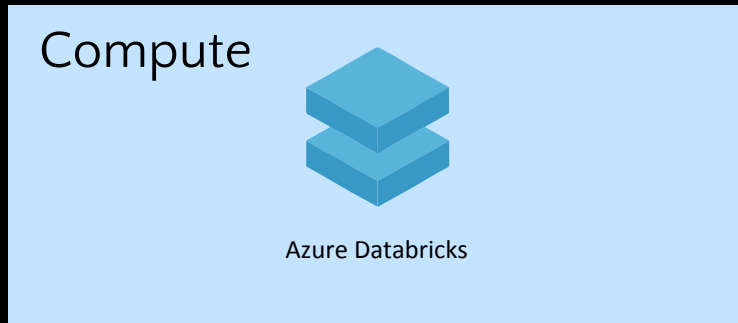# Azure Databricks Security | Data Protection

**Data Protection**

- Encryption-At-Rest – Service Managed Keys, User Managed Keys
- Encryption-in-flight (Transport Layer Security TLS)
- File/Folder Level ACLs for AAD Users, Groups, Service Principals
- ACLs for Clusters, Folders, Notebooks, Tables, Jobs
- Secrets with Azure Key Vault

* Preview feature

databricks

# Data Protection | Encryption | At-Rest

- Azure Databricks has separation of compute and storage

Compute

Azure Databricks

Storage

Azure BLOBs     Data Lake

- Storage Services such as Azure Blob Store, Azure Data Lake Storage Provide
  - Encryption of Data
  - Customer Managed Keys
  - File/Folder Level ACLs (Azure Data Lake Storage)

databricks

# Data Protection | Encryption | In-Transit

All the traffic from the Control Plane to the Clusters in the customer subscription is always encrypted with TLS.

# Data Protection | Access Control | ADLS Passthru

- Authenticate automatically to Azure Data Lake Storage (ADLS) from Azure Databricks clusters using the same Azure Active Directory (Azure AD) identity that one uses to log into Azure Databricks.

- Commands running on a configured cluster will be able to read and write data in Azure Data Lake Storage without requiring one to configure service principal credentials.

Azure Data Lake Storage Gen1 Credential Passthrough ❓

☑ Enable credential passthrough and only allow Python and SQL commands

# Data Protection | Access Control | Folders

| Ability | No Permissions | Read | Run | Edit | Manage |
|---|---|---|---|---|---|
| View items | | X | X | X | X |
| Create, clone, import, export items | | X | X | X | X |
| Run commands on notebooks | | | X | X | X |
| Attach/detach notebooks | | | X | X | X |
| Delete items | | | | X | X |
| Move/rename items | | | | X | X |
| Change permissions | | | | | X |

databricks

# Data Protection | Access Control | Notebooks

| Ability | No Permissions | Read | Run | Edit | Manage |
|---|---|---|---|---|---|
| View cells | | X | X | X | X |
| Comment | | X | X | X | X |
| Run commands | | | X | X | X |
| Attach/detach notebooks | | | X | X | X |
| Edit cells | | | | X | X |
| Change permissions | | | | | X |

databricks®

# Data Protection | Access Control | Clusters

| Ability | No Permissions | Can Attach To | Can Restart | Can Manage |
|---|---|---|---|---|
| Attach notebook to cluster | | x | x | x |
| View Spark UI | | x | x | x |
| View cluster metrics | | x | x | x |
| Terminate cluster | | | x | x |
| Start cluster | | | x | x |
| Restart cluster | | | x | x |
| Edit cluster | | | | x |
| Attach library to cluster | | | | x |
| Resize cluster | | | | x |
| Modify permissions | | | | x |

databricks

# Data Protection | Access Control | Jobs

| Ability | No Permissions | Can View | Can Manage Run | Is Owner | Can Manage (admin) |
|---|---|---|---|---|---|
| View job details and settings | X | X | X | X | X |
| View results, Spark UI, logs of a job run | | X | X | X | X |
| Run now | | | X | X | X |
| Cancel run | | | X | X | X |
| Edit job settings | | | | X | X |
| Modify permissions | | | | X | X |

databricks

# Data Protection | Access Control | Tables

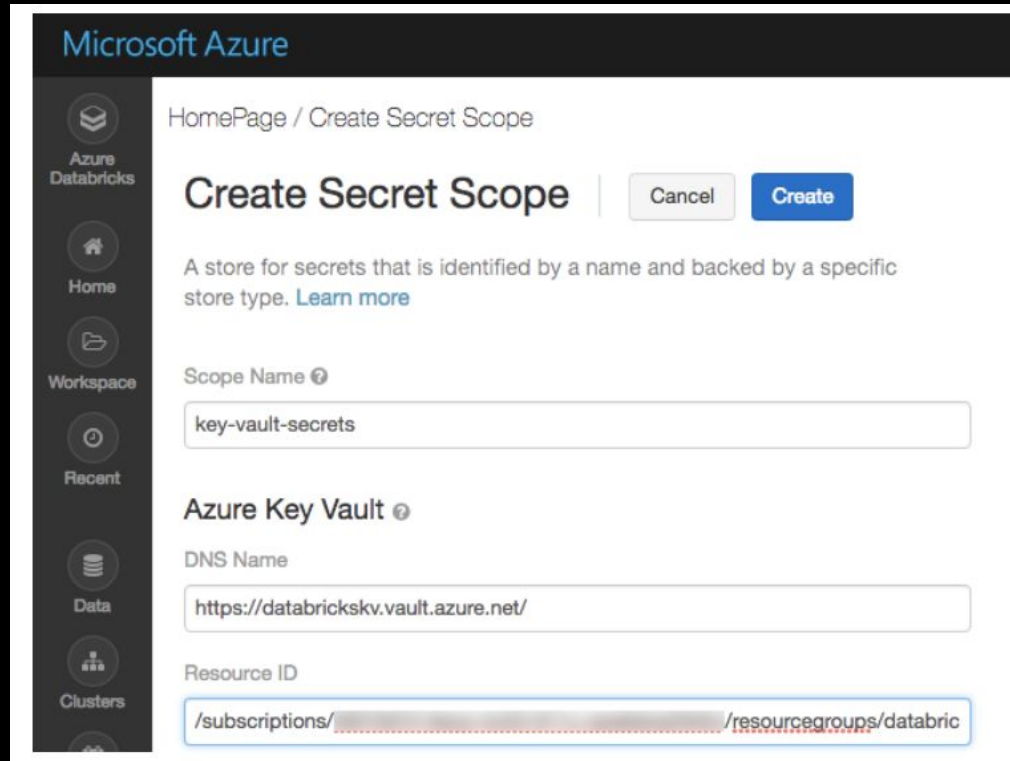**Objects**

CATALOG | DATABASE | TABLE | VIEW | FUNCTION | ANONYMOUS
FUNCTION | ANY FILE

**Privileges**

SELECT           – read access to an object
CREATE           – ability to create an object (eg. Table in a Database)
MODIFY           – ability to add/delete/modify data in an Object
READ_METADATA      –  ability to read Metadata about an object
ALL_PRIVELEGES      – all of the above

# Data Protection | Secrets

- Using our Secrets APIs, Secrets can be securely stored including in a Azure Key Vault or Databricks backend
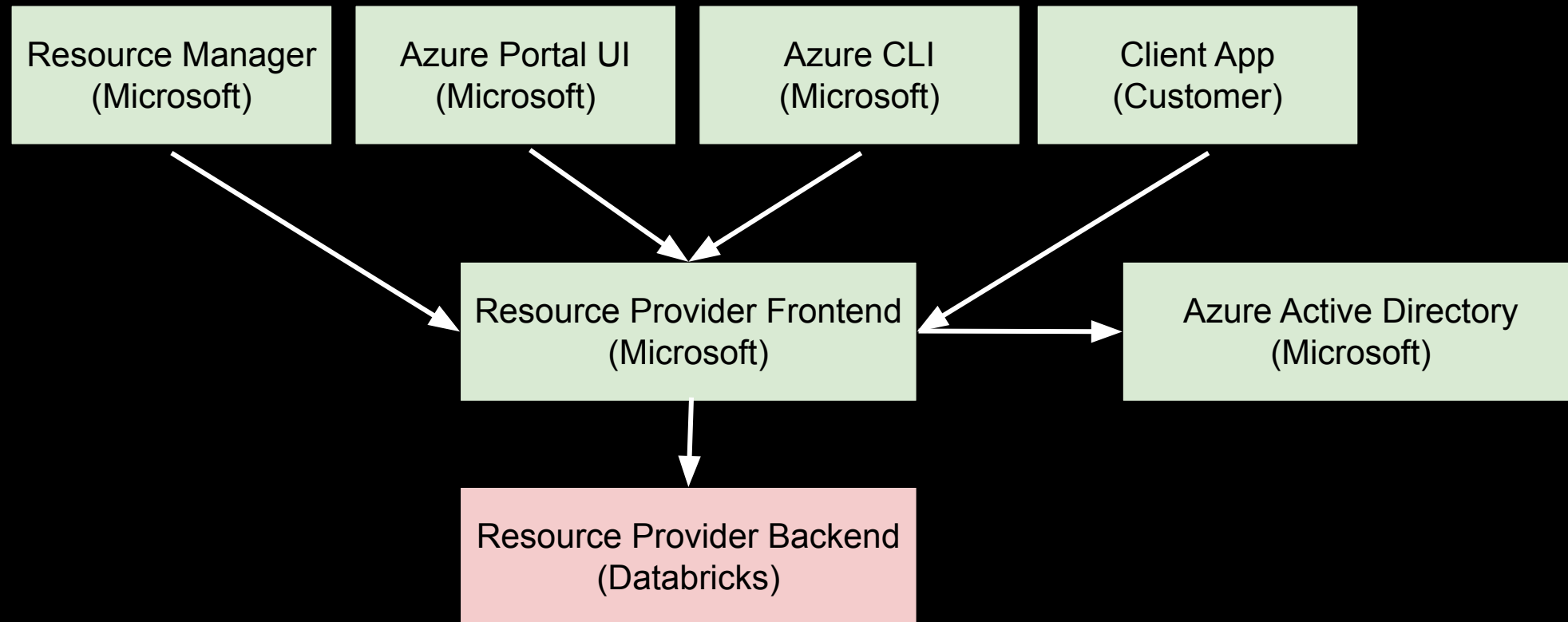- Authorized users can consume the secrets to access services

# Azure Databricks Security | IAM/Auth

IAM/Auth

- Azure Active Directory (AAD) Authentication (w/ MFA)
- AAD Groups (using SCIM)
- AAD Conditional Access
- AAD Access Tokens

databricks

# IAM/Auth | First-party AAD Integration

# IAM/Auth | SCIM Integration

Azure Databricks supports SCIM, or System for Cross-domain Identity Management, an open standard that allows you to automate user provisioning. SCIM lets you use Azure Active Directory to create users in Azure Databricks and give them the proper level of access, as well as remove access for users (deprovision them) when they leave the organization or no longer need access to Azure Databricks.



**Attribute Mappings**

Attribute mappings define how attributes are synchronized between Azure Active Directory and customappsso

| AZURE ACTIVE DIRECTORY ATTRIBUTE | CUSTOMAPPSS... | MATCHING ... | |
|---|---|---|---|
| userPrincipalName | userName | 1 | Delete |
| extensionAttribute1 | id | | Delete |
| mail | emails[type e... | | **Delete** |
| Join(" ", [givenName], [surname]) | displayName | | **Delete** |
| Switch([IsSoftDeleted], , "False", "True", "True", | active | | **Delete** |

Add New Mapping



**Attribute Mappings**

Attribute mappings define how attributes are synchronized between Azure Active Directory and customappsso

| AZURE ACTIVE DIRECTORY ATTRIBUTE | CUSTOMAPPSS... | MATCHING ... | |
|---|---|---|---|
| displayName | displayName | 1 | Delete |
| extensionAttribute1 | id | 2 | Delete |
| members | members | | **Delete** |

Add New Mapping

databricks

# IAM/Auth | Conditional Access

Azure Databricks supports Azure Active Directory conditional access, which allows administrators to control where and when users are permitted to sign in to Azure Databricks. For example, conditional access policies can restrict sign-in to your corporate network or can require multi-factor authentication.

# IAM/Auth | AAD Token Support

You could use AAD tokens to automate provisioning of Azure Databricks workspaces and access the Databricks REST API

**Private Preview Only**

# Azure Databricks Security | Network Security

Network Security

- Managed VNets
- VNet Peering
- VNET Injection/BYO VNET
  - On-Premises Data Access
  - Single-IP SNAT and Firewall-based filtering via custom routing
  - Service Endpoint

* Preview feature

databricks

# Network Security | VNET Peering

Virtual network (VNet) peering allows the virtual network in which your Azure Databricks resource is running to peer with another Azure virtual network. Traffic between virtual machines in the peered virtual networks is routed through the Microsoft backbone infrastructure, much like traffic is routed between virtual machines in the same virtual network, through private IP addresses only.

# Network Security | VNET Injection / BYO VNET

If you're looking to do specific network customizations, you could deploy Azure Databricks data plane resources in your own VNET

# Azure Databricks Platform Architecture

Deployment with VNET Injection and inter-node TLS (one could be used without the other)

Private, secure, encrypted with compute in your network

Plug into your identity system for seamless control

**Network Security**

**Identity and Access**

**E2**

**Compliance**

**Data Protection**

Battle tested for compliance and other sensitive data policies

Extend identity all the way down to the data natively

databricks

# Secure network connectivity



Databricks control plane services

New scalable relay enables large number of encrypted connections for 1000's of clusters

/cluster 1   /cluster 2   /cluster 3   /cluster 4

TLS 1.2 - Encrypted connections

Only outbound connections from enterprise environments.

Enterprise owned VPCs/VNETs

VPC/VNET

VPC/VNET

VPC/VNET

Network Security

databricks

**Enterprise Cloud Platform**

# Private Link

Databricks SaaS
control plane

Azure Network

Azure Private Link

All communication between control and data
plane is through the cloud providers private
network - never sent over a public network

NET

VNET

VNET

Enterprise owned VNETs running the data plane clusters

databricks

**Enterprise Cloud Platform**

# IP Access Lists

Control the networks that can access databricks

**HQ**

`216.58.76.12/28`

**BRANCH**

`72.12.84.112, 72.12.84.116`

`65.125.5.6`

Databricks control plane

UI

</API>

| ALLOW | DENY |
|---|---|
| 216.58.76.12/28 | 216.58.76.15 |
| 72.12.84.112 | |
| 72.12.84.116 | |

Dynamically update the list of networks

Deny IPs within a permitted subnet to shrink exposure

**Enterprise Cloud Platform**

databricks

# Bring your own VNET

Today

Q4

RDS
(Notebooks)

Peering

RDS
(Encrypted
Notebooks))

Existing enterprise VPC/VNET

New VPC/VNET

Existing
VPC/VNET

Workspace 2

Databricks managed clusters in
enterprise managed VPC/VNET

Connectivity to other applications
through peering

Local connectivity to other
applications without peering

# Data Centric Security

**Example:** End to end security for your data lake, defined in your metastore

Data Protection

**Identity Provider**

Validate permissions

**Access Rules**

**Hive Metastore**

Users only access views they are allowed to access

**Databricks**

Identity federation

Language of choice

python   R   SQL   scala

**Blob Storage**

databricks

**Enterprise Cloud Platform**

# Encryption with enterprise root keys

Data Protection

Customer Cloud

AWS KMS

Databricks control plane

Notebooks

Key Manager

Jobs

Get full control over keys used to encrypt data in the control plane. Revoking key revokes data access

Key hierarchy enables use of different keys for different notebooks

An audit log of key operations makes for easy reporting and audit requests

Customer creates key in Key Vault for an account or a workspace

Databricks creates data encrypting keys rooted in customer key

Applications use customer key to encrypt all the data

# End-to-end encrypted clusters

Data Protection

Databricks control plane

HTTPS

Cluster node

Proxy

Inter worker traffic = RPC + SSL (AES-128)

HTTPS

APACHE Spark™

Local disks

b3Blbm…za…1rZ…kto…EAAA AABG5…l…s…UA…A…bm9u ZQAAAAAAAABAAAABFw AAAAd2c1gtcnN…iAAAAw… EAA

Cache          Data

Enable encryption for a cluster - no key/cert management necessary

Ensure that data shared between cluster nodes is always encrypted

Guarantee that data is always encrypted no matter where its stored on the cluster - root, ephemeral or network attached disks

databricks

**Enterprise Cloud Platform**

# Security and Privacy

## Security

1. ISO 27001
2. SOC 2 Type 2
3. HITRUST (end of 2019 on Azure)
4. GovCloud & FedRAMP (2020 on Azure)

## Privacy

1. Privacy Shield
2. ISO 27018
3. GDPR compliant

**Enterprise Cloud Platform**

databricks

# Azure Databricks Security

**Compliance**



- Audit Logs
- ISO 27001
- ISO 27018
- HIPAA (Covered by MSFT BAA)
- SOC2, Type 2

databricks

# Compliance | Audit Logs

Databricks provides comprehensive end-to-end audit logs of activities performed by Databricks users, allowing your enterprise to monitor detailed Databricks usage patterns.

**Integration with Azure Monitor**

Services / Entities included are:

- Accounts
- Clusters
- DBFS
- Genie
- Jobs
- ACLs
- SSH
- Tables

# Reference Architectures

# Recommended End-to-End Architecture



Ingest     Store     Prep and train     Serve

Streaming data

Azure Event Hubs

Azure Machine Learning service

Model Serving

Azure Kubernetes service

Apps

Azure Databricks

Operational Databases

Cosmos DB, SQL DB

Ad–hoc Analysis

Power BI

Batch data

Azure Data Factory

Azure Data Lake Storage

Azure SQL data warehouse

Azure analysis services

databricks

# Recommended End-to-End Architecture - Batch ETL

# Recommended End-to-End Architecture - Stream ETL

# Recommended End-to-End Architecture - ML/Prediction

# Recommended End-to-End Architecture - BI and Analysis

# Recommended End-to-End Architecture



Ingest

Store

Prep and train

Serve

Streaming data

Azure Event Hubs

Batch data

Azure Data Factory

Azure Data Lake Storage

Azure Machine Learning service

Azure Databricks

Model Serving

Azure Kubernetes service

Operational Databases

Cosmos DB, SQL DB

Ad-hoc Analysis

Apps

Power BI

Azure SQL data warehouse

Azure analysis services

databricks

# Azure Databricks Best Practices

databricks

# Workspace Admin Best Practices

- Create different workspaces by different department / business team / data tier, and per environment (dev, qa, prod) - across relevant Azure subscriptions
- Define workspace level tags which propagate to initially provisioned resources in managed resource group (Tags could also propagate from parent resource group)
- Use ARM templates (search "databricks") to have a more managed way of deploying the workspaces - whether via CLI, powershell or some SDK
- Create relevant groups of users - using Group REST API or by using AAD Group Sync with SCIM

databricks

# Security Best Practices

- Do not store any production data on DBFS (use it only for toy / experimental datasets).
- Configure encryption-at-rest for Blob Storage and ADLS, preferably by using customer-managed keys in Azure Key Vault.
- Use Secrets with Azure Key Vault backend to obfuscate passwords and keys in notebooks.
- Prefer to use ADLS credential passthrough over Table ACLs (if possible).
- Configure access control for Databricks-native resources (clusters, notebooks, jobs etc.)
- Deploy workspace in your VNET to enable networking customizations.
- Configure Audit Logs to monitor the activity in a workspace.

# Tools & Integration Best Practices

- Use Azure Data Factory to orchestrate pipelines / workflows (or something like Airflow).
- Connect your IDE or custom applications to Azure Databricks clusters using DB-Connect (Private Preview).
- Sync notebooks with Azure Devops for seamless version control.
- Use Databricks CLI for CI / CD from relevant enterprise tools/products, or to integrate with other systems like on-prem SCM or Library Repos etc.
- Use Library Utilities to install python libraries scoped at notebook level (cluster-scoped libraries may make more sense in certain cases).
- Use Init Scripts to do custom installs at cluster level.

# Databricks Runtime Best Practices

- Use Delta wherever you can, to get the best performance and reliability for your big data workloads, and to create no-fuss multi-step data pipelines.
- Use Machine Learning Runtime for working with the latest ML/DL libraries (including HorovodRunner for distributed DL).
- Use DBIO Cache for accelerating reads from Blob Storage or ADLS.
- Use ABS-AQS connector for structured streaming when working with consistent rate of incoming files on Blob Storage.
- Turn on Databricks Advisor for automated tips on how to optimize workload processing.

# HA and DR Best Practices

- Deploy Azure Databricks in two paired azure regions, ideally mapped to different control plane regions.
  - E.g. East US2 and West US2 will map to different control planes
  - Whereas West and North Europe will map to same control plane
- Use Azure Traffic Manager to load balance and distribute API requests between two deployments, when the platform is primarily being used in a backend non-interactive mode.
- Design to honor API and other limits of the platform.
  - Max API calls/ hr = 1500
  - Jobs per hour per workspace = 1000
  - Maximum concurrent Notebooks per cluster = 145

# Cluster Best Practices

- Use autoscaling and auto-termination wherever applicable (e.g. auto-termination doesn't make sense if you need a cluster for data analysis by multiple users almost through the day, etc.).
- Use latest Databricks Runtime version to take advantage of latest performance & other optimizations (applicable in most cases, though not all).
- Use High-concurrency cluster mode for data analysis by a team of users via notebooks or a BI tool, or if you want to enforce data protection via Table ACLs or ADLS Passthrough.
- Use cluster tags for project / team based chargeback.

databricks

# Cluster Best Practices Contd..

- Use Spark config tab if certain tuning would make sense for a specific workload (like config to use broadcast join).
- Use Event Log and Spark UI to see how different queries / workload executions perform, and what affect those have on a cluster's health.
- Configure Cluster Log Delivery
- Use Cluster ACLs to configure what each user or a group of users are allowed to do.
- Refer this blog by a customer, which more or less mentions what we've covered here. Rest is really workload dependent where it requires evidence-based tuning.

# Appendix - Choosing the instance type

# Different Azure Instance Types

## Compute Optimized

- Fs
  - Haswell processor (Skylake not supported yet)
  - 1 core ~ 2GB RAM
  - SSD Storage: 1 core ~ 16GB
- H
  - High-performance
  - 1 core ~ 7GB RAM
  - SSD Storage: 1 core ~ 125GB

## Memory Optimized

- DSv2
  - Haswell processor
  - 1 core ~ 7GB RAM
  - SSD Storage: 1 core ~ 14 GB
- ESv3
  - High-performance (Broadwell processor)
  - 1 core ~ 8GB RAM
  - SSD Storage: 1 core ~ 16GB

## Storage Optimized

- L
  - 1 core ~ 8GB RAM
  - SSD Storage: 1 core ~ 170GB
  - Price : .156

## General Purpose

- DSv2 and DSv3
  - DSv2 - 1 core ~ 3.5GB RAM
  - DSv3 - 1 core ~ 4GB RAM
  - SSD Storage:
    - DSv2 - 1 core ~ 7GB
    - DSv3 - 1 core ~ 8GB

databricks

# Cluster Sizing Starting Points

## Rules of Thumb

- Fewer big instances > more small instances
  - Reduce network shuffle; Databricks has 1 executor / machine
  - Applies to batch ETL mainly (for streaming, one could start with smaller instances depending on complexity of transformation)
  - Not set in stone, and reverse would make sense in many cases - so sizing exercise matters
- Size based on the number of tasks initially, tweak later
  - Run the job with a small cluster to get idea of # of tasks (use 2-3x tasks per core for base sizing)
- Choose based on workload (Probably start with F-series or DSv2):
  - ETL with full file scans and no data reuse - F / DSv2
  - ML workload with data caching - DSv2 / F
  - Data Analysis - L
  - Streaming - F

# How do we tweak these?

**Workload requires caching (like machine learning)**

- Look at the Storage tab in Spark UI to see if the entirety of the training dataset is cached
  - Fully cached with room to spare -> less instances
  - Partially cached
    - Almost completely cached? -> Increase the cluster size
    - Not even close to cached -> Consider L series or DSv2 memory-optimized
      - Check to see if persist is MEMORY_ONLY, or MEMORY_AND_DISK
      - Spill to disk with SSD isn't so bad
- Still not good enough? Follow the steps in the next section

# How do we tweak these?

**ETL and Analytic Workloads**

- Are we compute bound?
  - Check CPU Usage (Ganglia metrics to come to Azure Databricks soon)
  - Only way to make faster is more cores

- Are we network bound?
  - Check for high spikes before compute heavy steps
  - Use bigger/fewer machines to reduce the shuffle
  - Use an ssd backed instance for faster remote reads

- Are we spilling a ton?
  - Check Spark SQL tab for spill (pre-agg before shuffles are common to spill)
    - Use L-series
    - Or use more memory

# Q&A