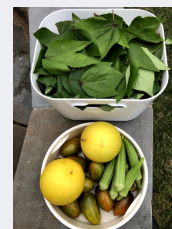# whoarewe

- David Veuve
  - Head of Security Field Engineering
- At work I am a
  - Security nerd
  - Content creator
  - Enabler
- I help with
  - Databricks platform security
  - Security detection and response
- I survived quarantine thanks to

- Arun Pamulapati
  - Sr. Staff Security Field Engineer
- At work I am a
  - Friends of the field
  - Builder
  - Collaborator
- I help with
  - Databricks platform security
  - Security use cases (Okta, DNS ...log analytics.)
- I survived quarantine thanks to

# Thesis

Databricks has built the security features
Databricks has helped 1000s of customers
Databricks has seen what works

Let's just tell you what works

35 minutes
to cover:

- 60+ Slides
- 17 Page Whitepaper
- Self-service tool

# What you will learn today

- Databricks Lakehouse Architecture
- Top threats impacting Databricks
- Example controls
- Where to get the full whitepaper
- How to analyze an existing deployment

# tl;dr

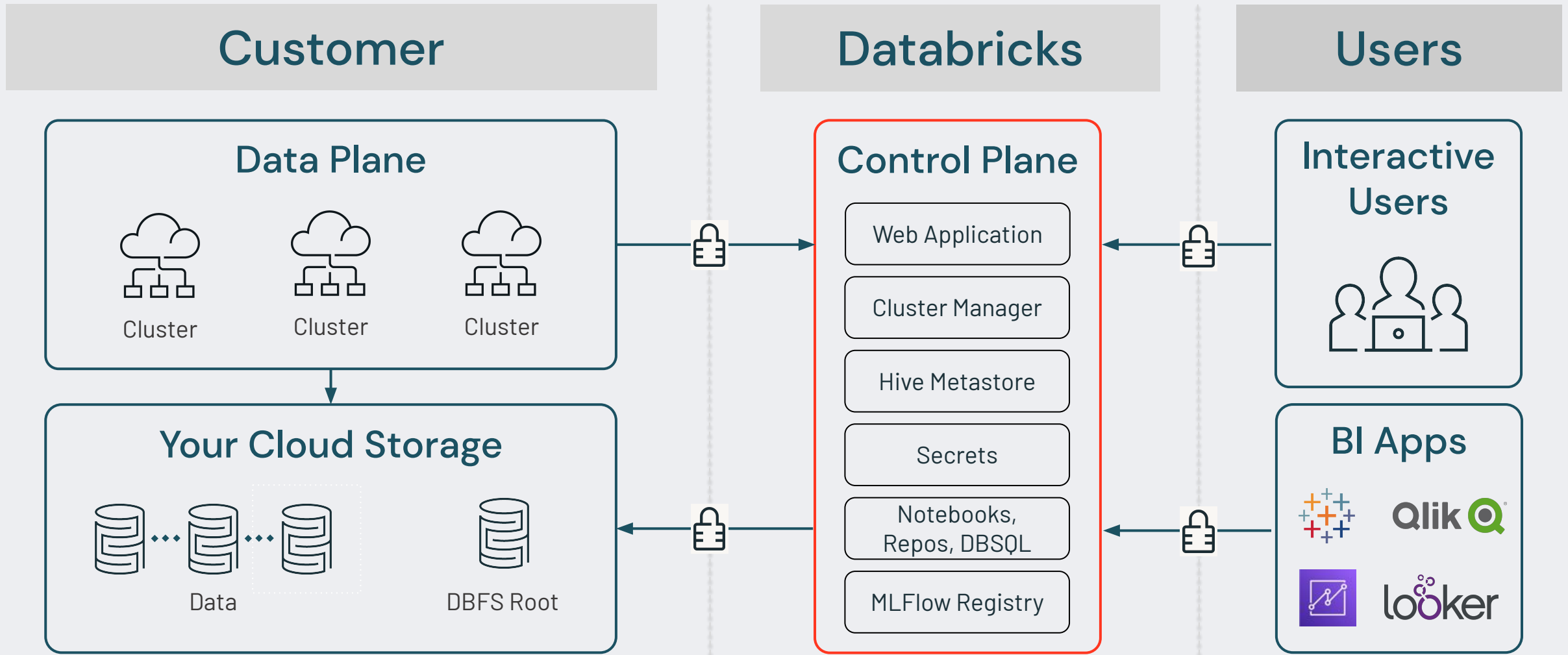| | |
|---|---|
| ## Security Architecture and Controls | ## Prescriptive Best Practices Guides |
| The security & trust center outlines our architecture and enumerates the key security features.<br><br>Overview: databricks.com/trust/whitepaper<br>Detail: databricks.com/trust | Databricks has captured our best practices into a doc with "common" and "high-security" models and checkboxes. *Check those boxes*.<br><br>databricks.com/trust/Security-Best-Practices_Databricks-on-AWS.pdf |
| ## Meet Compliance Needs | ## Analyze your deployment |
| Have your security team download our due diligence package (ISO certs, Pen Test) and reach out to your account team for SOC 2, Enterprise Security Guide<br><br>Compliance: databricks.com/trust#compliance | The Workspace Analysis tool now contains a security section! Check your workspace against our most common best practices.<br>Demo at the end of this talk!<br><br>Coming soon to the Databricks Blog! |

# Databricks Lakehouse Architecture

# What Are The Threats?

# Account Takeover

## Attackers gain the credentials/access of your users
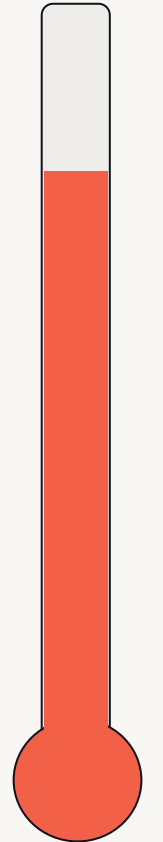
### Risk Overview

- Customers often analyze sensitive datasets
- Compromised end-user credentials grant access
    - Phishing, brute force, etc.

### Best Practices to Mitigate

- Required two-factor auth on your identity provider
    - Consider FIDO Key
- Manage local passwords
- Use SCIM to deprovision
- IP Access Lists or PrivateLink
- Monitor Audit Log
- Limit token lifetime

Attack Likelihood

# Accidental Insider Exposure

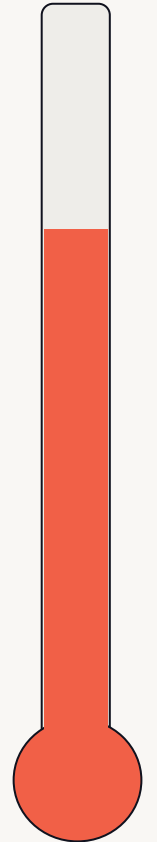## Resource abuse or accidental copy by insiders

### Risk Overview

- Someone believes their job is easier without pesky security controls
- Data is copied where it shouldn't be, or where ACLs aren't applied

### Best Practices to Mitigate

- Backup your data and code
- Run most sensitive ETL through a CI/CD process (code review)
- Utilize secured models like Table ACLs for limited exposure
- Limit data in DBFS and monitor for large datasets
- Deploy data exfiltration protections

Attack Likelihood

# Data Exfiltration

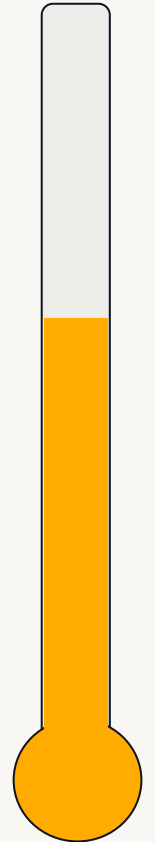## Data stolen by an attacker or malicious insider

### Risk Overview

– Risk is of a user sending sensitive data out to some external location

– Data plane needs connectivity, by default has full outbound

– Also consider access to cloud storage (S3 buckets, etc.)

### Best Practices to Mitigate

– Use Customer-managed VPC / VNet Injection (baseline)

– Route traffic through a firewall or proxy to limit destination

– [AWS] VPC Endpoint Policies

– Limit access to sensitive data

– Configure data exfiltration settings in the console

# Also covered in the docs

## Resource Abuse

- Customer cloud infrastructure hijacking for crypto mining
- Accidental/abusive waste of customer resources

## Compromise of Databricks, Inc.

- Compromise of Databricks Inc user or system could result in compromise of customer environment

# What is a security team to do?

# We want to securely authenticate users

# Authenticate via single sign-on

DE

AIE

RA

CD

Databricks lakehouse platform

Data Science and ML

Real-Time Data Applications

...nt and Governance

...Data Lake

...ty & Administration

Users

## Use Multi-factor Authentication as well!

- SAMLv2, SCIM support
- Documented integrations with six partners

G

AWS

# SCIM & Role-based Access Control



Admins set ACLs via the UI or Permissions API

Account Takeover

DE

Accidental Insider Exposure

RA

CD

Databricks Control Plane

VIEW

CREATE

MANAGE

ATTACH

RUN

DELETE

Sync using SCIM

Notebooks

Jobs

Clusters

ACLs for:
- Jobs
- Notebooks
- Clusters
- Pools
- Tables
- Workspaces
- Secrets
- ML Experiments

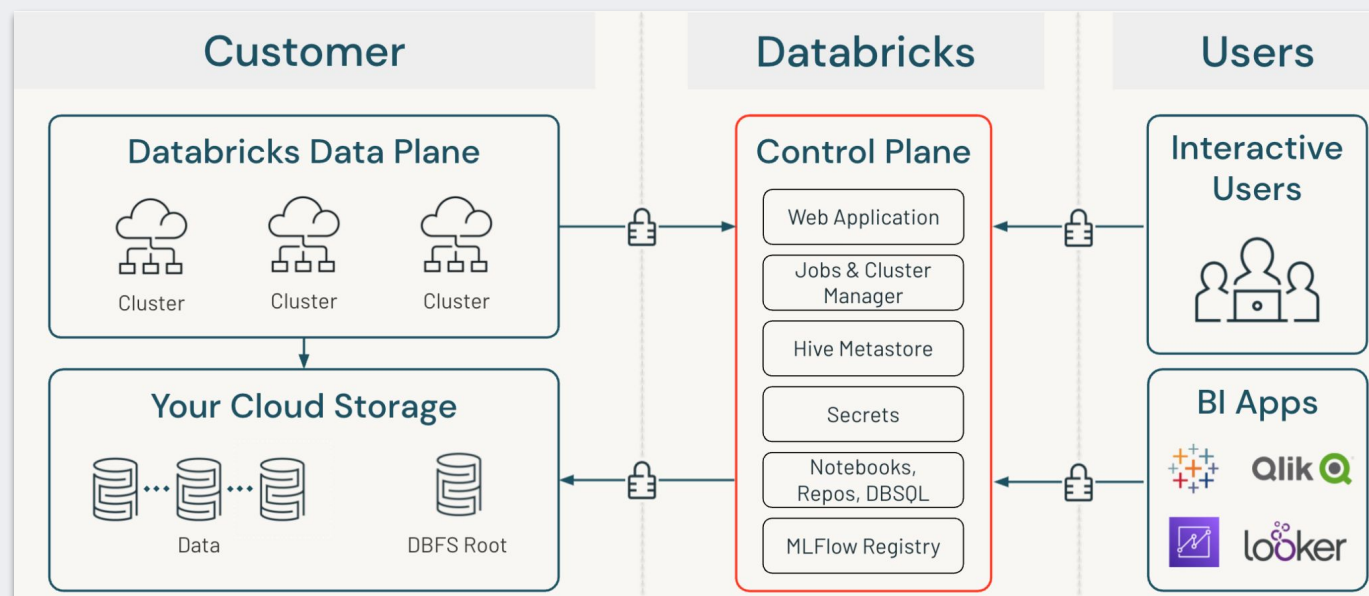# We want to encrypt data

# Encrypt buckets and restrict access

- Encrypt your buckets
- Don't allow the public to access your buckets

# Customer-managed keys
## (AWS, Azure)

**Compromise of Databricks Inc.**

Customers can enable customer-managed keys

- CMK for Managed Services

  Priority control plane data

- CMK for Workspace Storage

  DBFS root storage and (AWS-only) EBS volumes



Customer Cloud

KMS / AKV

Databricks control plane

Development / Experimentation

Key Manager

Production Jobs

Customer creates key in Cloud-native Key Management Service for an account or a workspace

Databricks data plane

DBFS Storage

Cluster Volumes

# We want to isolate different workloads

# Design workload isolation

## Workspace 1

Cluster 1

Cluster 2
(Table ACLs)

## Workspace 2

Cluster 3

Easily isolate based on your needs

1. Table ACLs supports users with different privileges
2. Multiple workspaces isolate groups who won't collaborate
3. Limit usage of standard clusters

# Avoid storing production data in DBFS

- DBFS (Databricks Filesystem) is accessible to all workspace users
- Instead, store data in buckets / ADLS
- AWS customers can use bucket policy to limit access



Databricks Data Plane

Cluster     Cluster     Cluster

Your Cloud Storage

Data                    DBFS Root

# We want to prevent data loss

# Deploy with a customer-managed VPC

## Azure name: VNet Injection

- Limit outgoing connections: Use a firewall or proxy to limit outbound traffic.
- (AWS-only) Lower privilege level: maintain more control of your own AWS account.
- Simplified network operations: Better network space utilization.
- Consolidation of VPCs: Multiple data planes can share a single VPC

# Implement data loss protections

- Lock down outbound access with customer managed VPC
- Restrict outbound access
- Restrict access to storage

# We want to have a record of what happens with our data and detect user compromise

DATA+AI
SUMMIT 2022

# Databricks audit log

Account Takeover

Data Exfiltration

Accidental Insider Exposure

Resource Abuse

Compromise of Databricks Inc.

## Audit Logging

- Customers can configure near-real-time logging
  - (AWS/GCP) to a bucket owned by the customer
  - (Azure) to diagnostic logging
- (AWS) Cloudtrail logs also includes provisioning activities

## System Logs

- Understand system activities via system logs, including stdout, stderr, etc.
- Use metrics to understand utilization and health

Ready-to-use analysis notebooks
on our blog!
(Linked from whitepaper)

# We want to avoid loss of data

**DATA+AI**
SUMMIT 2022

# Backup your control plane data

- Backup via the Databricks migration tool
- Databricks Repos moves your code storage directly to git

**DATA+AI**
SUMMIT 2022

# Manage code run via CI/CD

Mature organizations will often build their production workloads via CI/CD:

- [CI/CD is very compatible](#) in a Databricks environment.
- For security:
  - Integrate code scanning, better provide for permissions, perform linting, and more
  - Scan for passwords, token in the code

# We want to avoid leaking passwords

# Store and use secrets securely

- Securely store API Keys or other credentials and reference by location
- Secrets help avoid hard-coded credentials in code, but doesn't hide them from authorized users

```
api_key = dbutils.secrets.get(
                scope="myScope",
                key="mySecret"
        )
```

# Use AWS Nitro instances

AWS Nitro instances provide two major security benefits:

- NVMe disks automatically encrypt at rest (local SSDs for shuffle data)
  - Azure can enable software encryption
  - GCP automatically has encryption
- Many Nitro instances automatically encrypt data in transit between hosts
  - Azure can enable intra-cluster encryption in software
  - GCP automatically has inter-host encryption

**DATA+AI**
SUMMIT 2022

# Control user network access

Two options:

- IP Access Lists  Enterprise
- User to Control Plane Private Link (Preview)  Enterprise
  - Not available on GCP

**DATA+AI**
SUMMIT 2022

# Configure the admin console settings

Manage key configuration settings

- Can Databricks staff access your workspace for support?
- Can notebooks be exported?
- Can results be downloaded?
- Where are query results stored?
- Are ACLs enabled?

Admin Console

Users   Groups   Instance Profiles   Workspace Storage   Access Control   Single Sign On   Advanced   Global Init Scripts

Workspace Access Control: **Disabled**
What this means ❯

Cluster, Pool and Jobs Access Control: **Disabled**
What this means ❯

Table Access Control: **Disabled**
What this means ❯

Customer Approved Workspace Login: **Disabled** [Enable]
What this means ❯

Personal Access Tokens: **Enabled** [Disable] [Permission Settings]
What this means ❯

# Token management

- Enable or disable personal access tokens (PATs) for some or all users
- Configure a max token lifetime for new tokens

# We want to utilize private networking

**DATA+AI**
SUMMIT 2022

# Configure Back-End Private Link

Multiple options:

1. Data Plane to Control Plane
   - By default, the cloud service provider backbone is very secure
   - (AWS / Azure) Use Data Plane to Control Plane PrivateLink (Preview) Enterprise
2. Data Plane to Storage Services
   - Public Endpoints
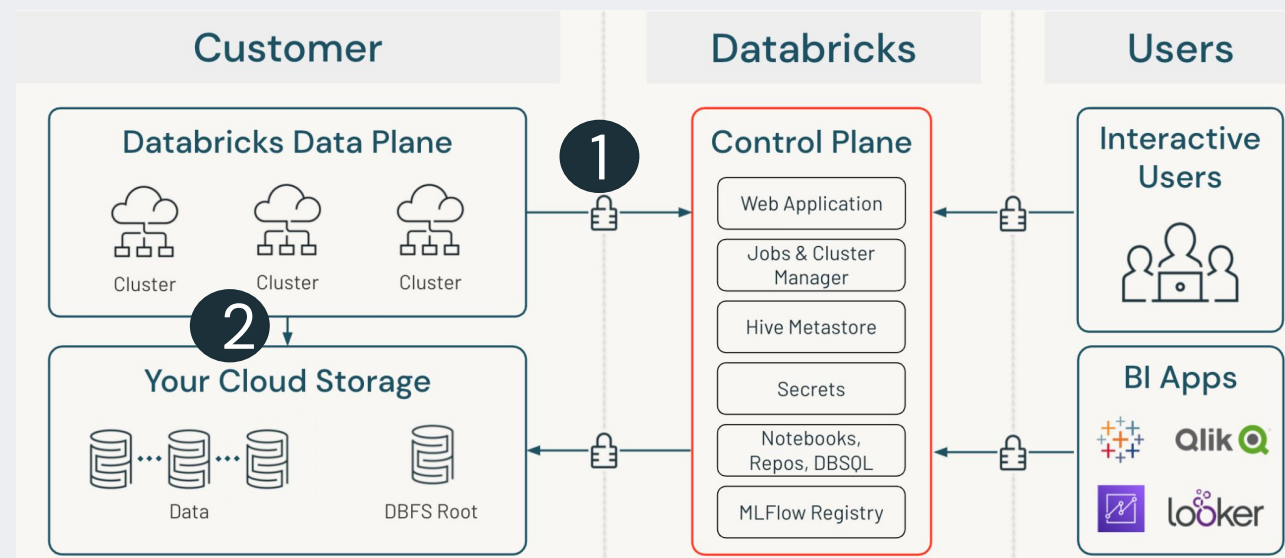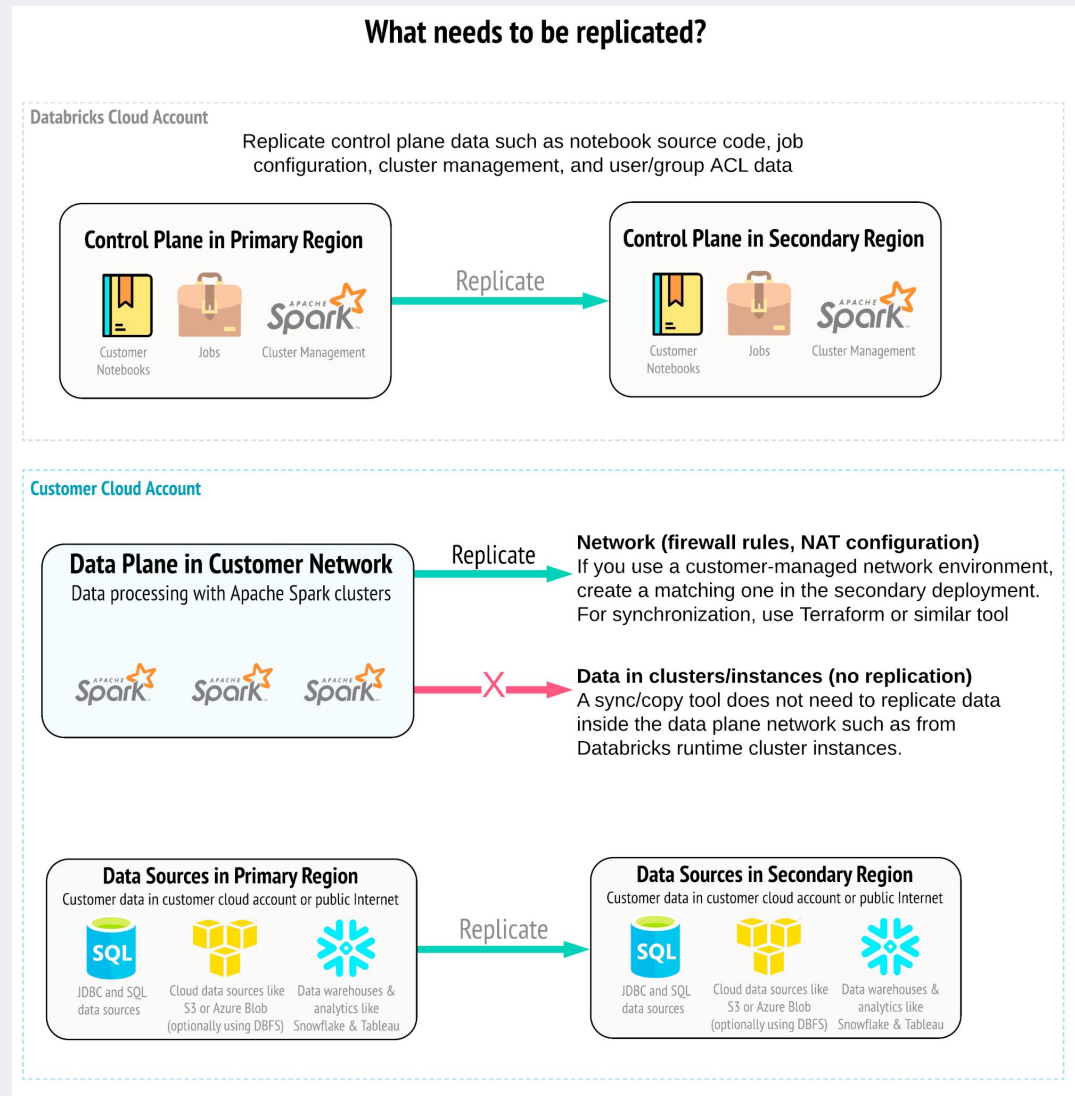   - Private Endpoints
   - Private Link to your data sources

# Configure a DR site

- Understand your business needs
- Choose a process that meets your business needs
- Prep workspaces and do a one-time copy
- Prepare your data sources
- Implement and test your solution



**What needs to be replicated?**

Databricks Cloud Account

Replicate control plane data such as notebook source code, job configuration, cluster management, and user/group ACL data

**Control Plane in Primary Region**
Customer Notebooks — Jobs — Cluster Management

Replicate →

**Control Plane in Secondary Region**
Customer Notebooks — Jobs — Cluster Management

Customer Cloud Account

**Data Plane in Customer Network**
Data processing with Apache Spark clusters

Replicate →

**Network (firewall rules, NAT configuration)**
If you use a customer-managed network environment, create a matching one in the secondary deployment. For synchronization, use Terraform or similar tool

X →

**Data in clusters/instances (no replication)**
A sync/copy tool does not need to replicate data inside the data plane network such as from Databricks runtime cluster instances.

**Data Sources in Primary Region**
Customer data in customer cloud account or public Internet
JDBC and SQL data sources — Cloud data sources like S3 or Azure Blob (optionally using DBFS) — Data warehouses & analytics like Snowflake & Tableau

Replicate →

**Data Sources in Secondary Region**
Customer data in customer cloud account or public Internet
JDBC and SQL data sources — Cloud data sources like S3 or Azure Blob (optionally using DBFS) — Data warehouses & analytics like Snowflake & Tableau

DATA+AI
SUMMIT 2022

# We want to control configurations and costs

# Configure cluster policies

- Limit users to create clusters with prescribed settings.
- Simplify the user interface for your users
- Control cost by limiting per cluster maximum cost

# Configure tagging to monitor

- Monitor cost and accurately attribute Databricks usage to your organization's business units and teams (for chargebacks, for example), you can tag clusters and pools.
- These tags propagate both to detailed DBU usage reports and to your cloud service provider (e.g., AWS EC2 and AWS EBS)

**DATA+AI**
SUMMIT 2022

# Security Best Practices Documentation

## Including a checklist! databricks.com/trust

Security Best Practices for Databricks on AWS

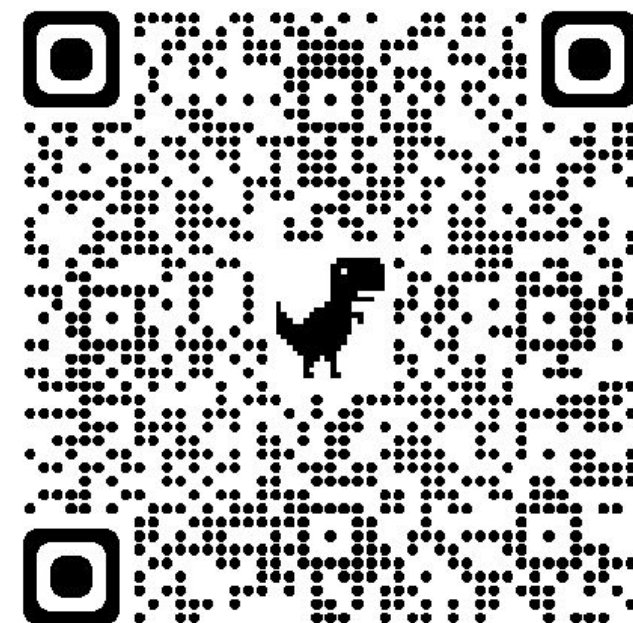Version 1.0 - June 16 2022

databricks

---

The following typical configurations are part of most enterprise production Databricks deployments. If you are a small data science team of a few people, you may not feel the need to deploy all of these. If Databricks may become a key pa of your business or if you are analyzing sensitive data, we recommend that you review these.

☐ Evaluate whether multiple workspaces are required for segmentation
☐ Check that your S3 buckets are encrypted and that public access is blocked
☐ Deploy Databricks into a customer-managed VPC for increased control over the network environment. Even if you do not need this now, this option increases the chances for future success with your initial workspace
☐ Authenticate via single sign-on
☐ Use multi-factor authentication
☐ Separate accounts with admin privileges from day-to-day user accounts
☐ Configure Databricks audit log delivery
☐ Configure maximum token lifetimes for future tokens using token management
☐ Configure admin console settings according to your organization's needs
☐ Apply bucket policies or other mitigations to avoid storing production datasets in DBFS
☐ Backup your notebooks stored in the control plane or store your notebooks in git repos
☐ Store and use secrets securely in Databricks or using a third-party service
☐ Consider whether to implement network protections for data loss

### Highly-secure deployments

In addition to the configurations typical to all deployments, the following configurations are often used in highly-secure Databricks deployments. While these are common, not all highly-secure environments use all of these settings. We recommend incorporating these items and the threat model in the following section alongside your existing security practices.

☐ Evaluate whether customer-managed encryption keys are needed on the control plane or data plane for contro over data at rest (Requires Enterprise tie)
☐ Keep an up-to-date user list by using SCIM
☐ Set complex local passwords or disable local passwords
☐ Use either IP access lists or front-end PrivateLink
☐ Configure back-end (data plane to control plane) PrivateLink connectivity
☐ Implement network protections for data exfiltration
☐ Evaluate whether your datasets require bucket versioning
☐ Evaluate whether your workflow requires using git repos or CI/CD
☐ Plan for and deploy a disaster recovery site if you have strong continuity requirements
☐ Consider requiring AWS Nitro instances that provide encryption for ephemeral storage at rest and between instances
☐ Encourage the use of clusters that support user isolation
☐ Configure cluster policies to enforce data access patterns and control costs
☐ Evaluate tagging to monitor and manage chargeback and cost control

---

AWS only for now...

But your deployment is probably fine, right?

Want to check anything real quick?

**DATA+AI**
SUMMIT 2022

# ANNOUNCING

## SWAT

### Security Workspace Analysis Tool!

# Security Workspace Analysis Tool

**(AWS) Beta starting shortly...**

– Compare workspace configurations against specific best practices

– Flag deviations and score your workspace over a period of time

– Get mitigation references

– Want to try this in your workspace?

Contact us at:
cybersecurity@databricks.com

Anindita Mahapatra   Ramdas Murali   Arun Pamulapati

# Recap

# tl;dr

## Security Architecture and Controls

The security & trust center outlines our architecture and enumerates the key security detections.

Overview: databricks.com/trust/whitepaper
Detail: databricks.com/trust

## Prescriptive Best Practices Guides

Databricks has captured our best practices into a doc with "common" and "high-security" models and checkboxes. *Check those boxes.*

databricks.com/trust/Security-Best-Practices_Databricks-on-AWS.pdf

## Meet Compliance Needs

Have your security team download our due diligence package (ISO certs, Pen Test) and reach out to your account team for SOC 2, Enterprise Security Guide

Compliance: databricks.com/trust#compliance

## Analyze your deployment

The Workspace Analysis tool now contains a security section! Check your workspace against our most common best practices.

Coming soon to the Databricks Blog!

# Thank you!

# ⭐ SWAT - Security Workspace Analysis Tool  SWAT

Share  Subscribe ⌄  🔄 Refresh  ⋮

| AccountId | Deployment Name | Tier | Region | WS Status |
|---|---|---|---|---|
| a2033dd6█████████████de27970 | db-████████gtz-2 | ENTERPRISE | us-east-1 | RUNNING |
| 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago |

| Databases | Tables | Groups | Users | Databricks Jobs | Notebooks |
|---|---|---|---|---|---|
| **2** | **5** | **2** | **3** | **0** | **320** |
| #Databases | #Tables | #Groups | #Users | #Databricks Jobs | #Notebooks |
| 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago |

## Workspace Settings

| chk_name | sub_category | score | severity | details |
|---|---|---|---|---|
| WS-1 | enableJobViewAcls | 0 | OK | Job Visibility Control |
| WS-2 | enforceClusterViewAcls | 0 | OK | Cluster Visibility Control |
| WS-3 | enforceWorkspaceViewAcls | 0 | OK | Workspace Visibility Control |
| WS-4 | enableProjectTypeInWorkspace | 1 | High | Enable Repos |
| WS-5 | enableResultsDownloading | 1 | High | Download button for notebook results |

🔄 4 days ago

# Network Security

- NPIP
- BYOVPC
- VPC Peering
- IP Access Lists

| High | Medium | Low |
|---|---|---|
| **0** | **2** | **0** |
| 🔄 4 days ago | 🔄 4 days ago | 🔄 4 days ago |

Share    Subscribe ⌄    ↻ Refresh

| DP-4 | S3 Encryption | 1 | High | Encrypt S3 buckets and restrict access |

↻ 4 days ago

# Compliance

- Cluster Policy
- Audit Logs
- Global Init Script
- Mounts

| High | Medium | Low |
|------|--------|-----|
| 2 | 0 | 0 |
| ↻ 4 days ago | ↻ 4 days ago | ↻ 4 days ago |

## Compliance Details

| chk_name | sub_category | score | severity | details |
|----------|--------------|-------|----------|---------|
| CP-10 | Instance Pool Custom Tag | 0 | OK | Didnt detect any Instance Pool Custom tags |
| CP-11 | Max concurrent runs | 0 | OK | All max concurrent runs < 5 |
| CP-12 | Global libraries | 0 | OK | No global libraries |
| CP-13 | User Privileges | 0 | OK | Controlled cluster create privileges |
| CP-14 | Log delivery configurations. | 1 | High | Audit Log Delivery has not been Enabled |
| CP-3 | AllPurpose Cluster Custom Tags | 0 | OK | All AllPurpose Clusters have custom tags |
| CP-4 | Job Cluster Custom Tag | 0 | OK | All Job Clusters have custom tags |
| CP-5 | AllPurpose Cluster Log Conf | 0 | OK | All AllPurpose Clusters have Log configuration enabled |
| CP-6 | Job Cluster Log Conf | 0 | OK | All Job clusters have Log configuration enabled |
| CP-7 | Managed Tables | 0 | OK | No managed tables |

↻ 4 days ago

2/3

# Announcing:
Security
Workspace
Analysis
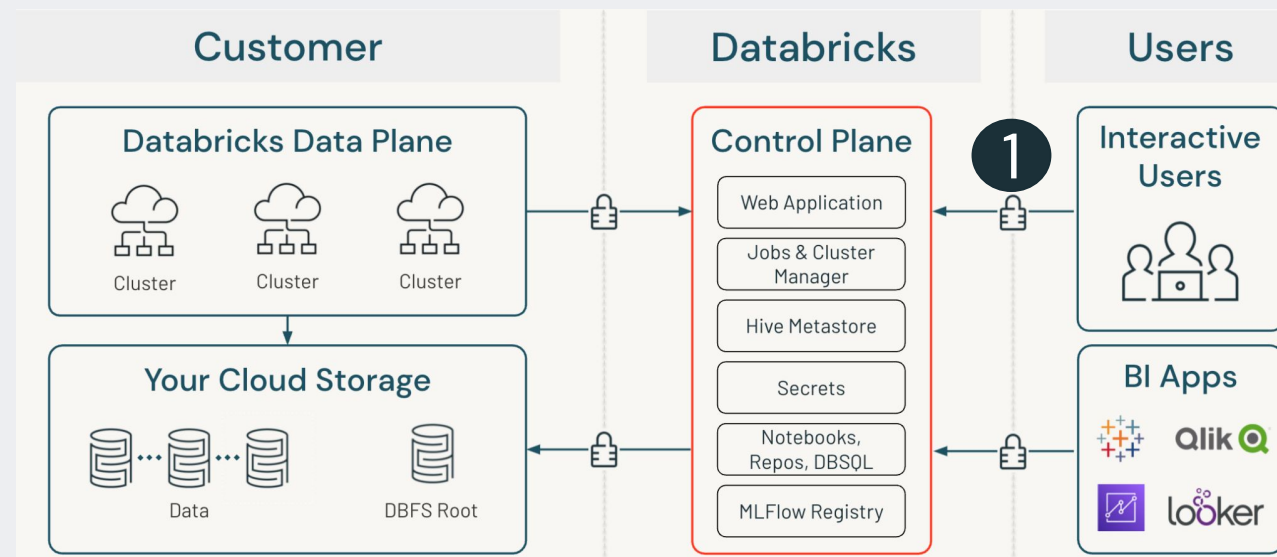Tool!

# Use AWS Nitro instances

AWS Nitro instances provide two major security benefits:

- NVMe disks automatically encrypt at rest (local SSDs for shuffle data)
  - Azure can enable software encryption
  - GCP automatically has encryption
- Many Nitro instances automatically encrypt data in transit between hosts
  - Azure can enable intra-cluster encryption in software
  - GCP automatically has inter-host encryption

# Control user network access

Two options:

- IP Access Lists  [Enterprise]
- User to Control Plane Private Link (Preview)  [Enterprise]
  - Not available on GCP

# Configure the admin console settings
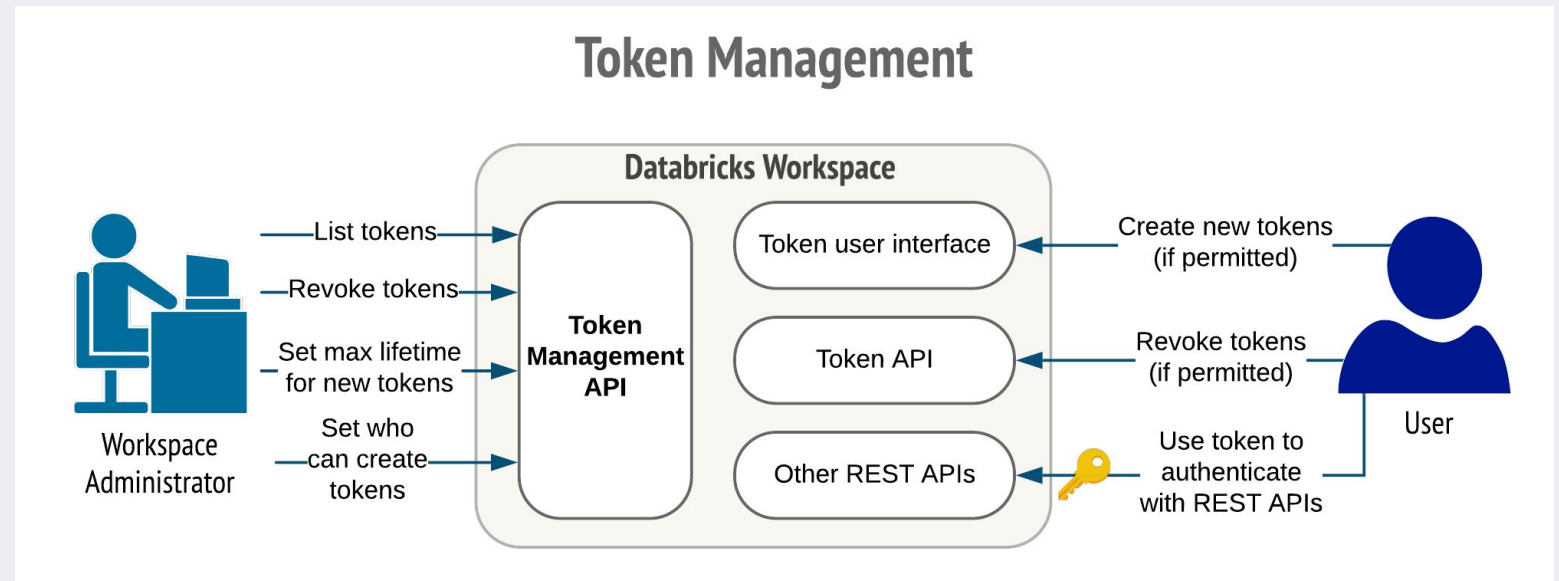
Manage key configuration settings

- Can Databricks staff access your workspace for support?
- Can notebooks be exported?
- Can results be downloaded?
- Where are query results stored?
- Are ACLs enabled?

Admin Console

Users  Groups  Instance Profiles  Workspace Storage  **Access Control**  Single Sign On  Advanced  Global Init Scripts

Workspace Access Control: **Disabled**
What this means ❯

Cluster, Pool and Jobs Access Control: **Disabled**
What this means ❯

Table Access Control: **Disabled**
What this means ❯

Customer Approved Workspace Login: **Disabled**  [Enable]
What this means ❯

Personal Access Tokens: **Enabled**  [Disable]  [Permission Settings]
What this means ❯

**DATA+AI**
SUMMIT 2022

# Token management

- Enable or disable personal access tokens (PATs) for some or all users
- Configure a max token lifetime for new tokens

# We want to utilize private networking
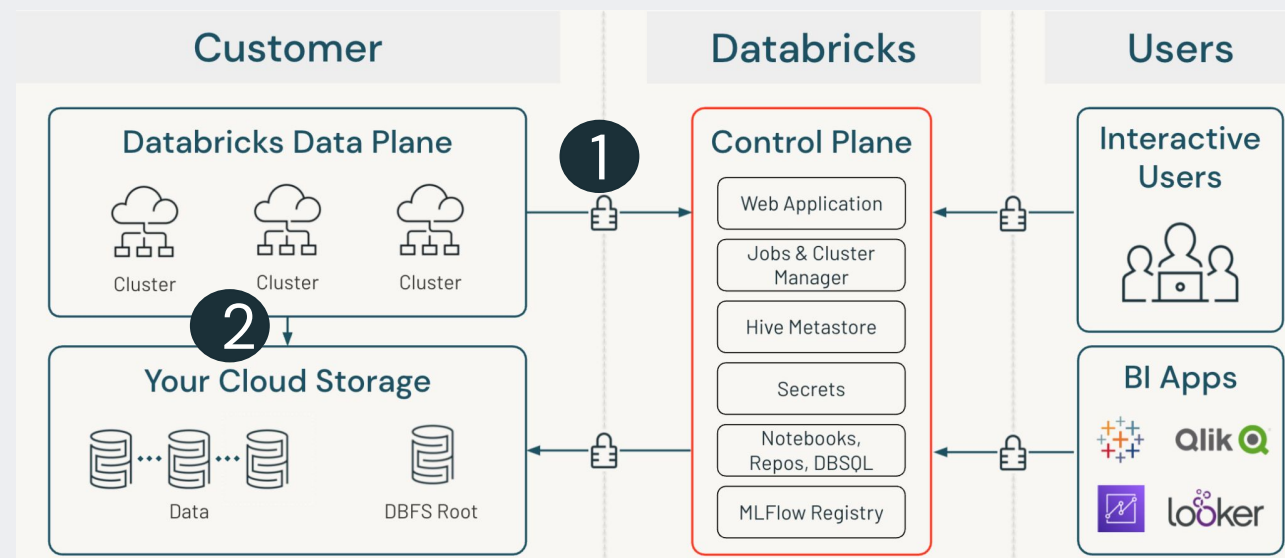
# Configure Back-End Private Link

Multiple options:

1. Data Plane to Control Plane
   - By default, the cloud service provider backbone is very secure
   - (AWS / Azure) Use Data Plane to Control Plane PrivateLink (Preview) [Enterprise]
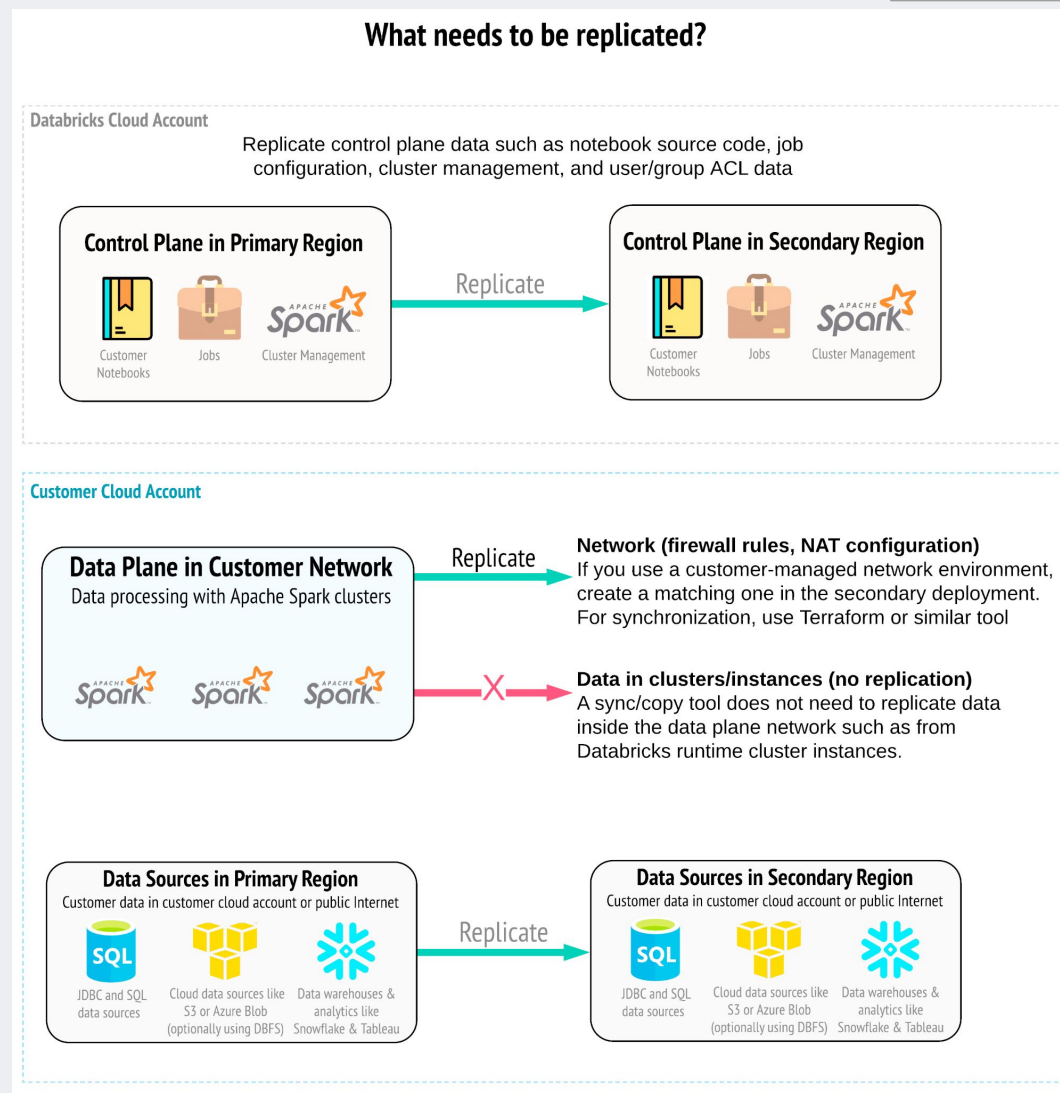2. Data Plane to Storage Services
   - Public Endpoints
   - Private Endpoints
   - Private Link to your data sources

# Configure a DR site

- Understand your business needs
- Choose a process that meets your business needs
- Prep workspaces and do a one-time copy
- Prepare your data sources
- Implement and test your solution



**What needs to be replicated?**

Databricks Cloud Account
Replicate control plane data such as notebook source code, job configuration, cluster management, and user/group ACL data

**Control Plane in Primary Region**
Customer Notebooks · Jobs · Cluster Management

Replicate →

**Control Plane in Secondary Region**
Customer Notebooks · Jobs · Cluster Management

Customer Cloud Account

**Data Plane in Customer Network**
Data processing with Apache Spark clusters

Replicate →
**Network (firewall rules, NAT configuration)**
If you use a customer-managed network environment, create a matching one in the secondary deployment. For synchronization, use Terraform or similar tool

X →
**Data in clusters/instances (no replication)**
A sync/copy tool does not need to replicate data inside the data plane network such as from Databricks runtime cluster instances.

**Data Sources in Primary Region**
Customer data in customer cloud account or public Internet
JDBC and SQL data sources · Cloud data sources like S3 or Azure Blob (optionally using DBFS) · Data warehouses & analytics like Snowflake & Tableau

Replicate →

**Data Sources in Secondary Region**
Customer data in customer cloud account or public Internet
JDBC and SQL data sources · Cloud data sources like S3 or Azure Blob (optionally using DBFS) · Data warehouses & analytics like Snowflake & Tableau

# We want to control configurations and costs

# Configure cluster policies

- Limit users to create clusters with prescribed settings.
- Simplify the user interface for your users
- Control cost by limiting per cluster maximum cost



Clusters / Cluster Policies / Create Policy

Create Cluster Policy    Cancel    Create

Name

Definition    Permissions

1

# Configure tagging to monitor

- Monitor cost and accurately attribute Databricks usage to your organization's business units and teams (for chargebacks, for example), you can tag clusters and pools.
- These tags propagate both to detailed DBU usage reports and to your cloud service provider (e.g., AWS EC2 and AWS EBS)

# Resource Abuse

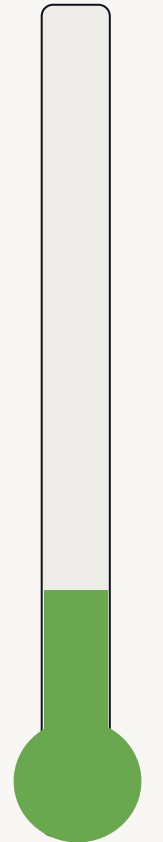`Attackers gain access to customer compute resources`

## Risk Overview

- Customer cloud infrastructure hijacking for crypto mining
- Accidental/abusive waste of customer resources

## Best Practices to Mitigate

- Cloud native protections:
    - Restricted x IAM role
    - Service quotas
    - Cloud monitoring
    - CloudTrail
- Databricks protections:
    - Cluster policies
    - Cluster ACLs
    - Library control
    - Databricks Audit logs

Attack Likelihood

# Compromise of Databricks Inc

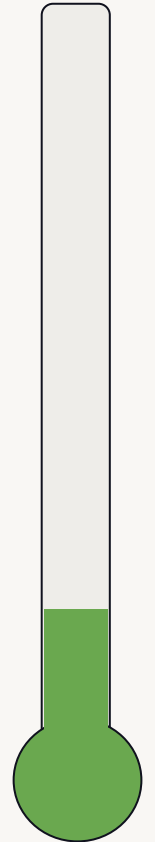Attacker gains customer environment via Databricks

Attack Likelihood

## Risk Overview

- Compromise of Databricks Inc user or system could result in compromise of customer environment

## Best Practices to Mitigate

- Monitor Databricks audit logs
- Consider CAWL
- Restrict cross account IAM role
- Monitor CloudTrail
- Break glass controls:
  - Customer managed key for managed services
  - Ability to disable Databricks cross account IAM role

# Secure local passwords
## (AWS)

- When MFA isn't available, use very long and complex passwords, and securely store them
- Disable local password via [password access control](#)

Permission Settings for: **Password Usage**   ×

NAME                    PERMISSION

admins                  Can Use   ▾   ✕

Select User, Group or Service Principal... ▾   Can Use ▾   + Add

Cancel   Save

Scripts   Workspace Settings

Permission Settings

ogin and API calls.

abled. Ensure single sign-on is configured correctly all users may be locked out of the workspace.

Permission Settings