# Databricks Serverless

**Nikhil Jethava**
Staff Product Manager, Databricks

**Matt Ryan**
Director Engineering and Co-founder, Kythera Labs

**Aaron Davidson**
Principal Software Engineer, Databricks

# Product safe harbor statement

This information is provided to outline Databricks' general product direction and is for informational purposes only. Customers who purchase Databricks services should make their purchase decisions relying solely upon services, features, and functions that are currently available. Unreleased features or functionality described in forward-looking statements are subject to change at Databricks discretion and may not be delivered as planned or at all.
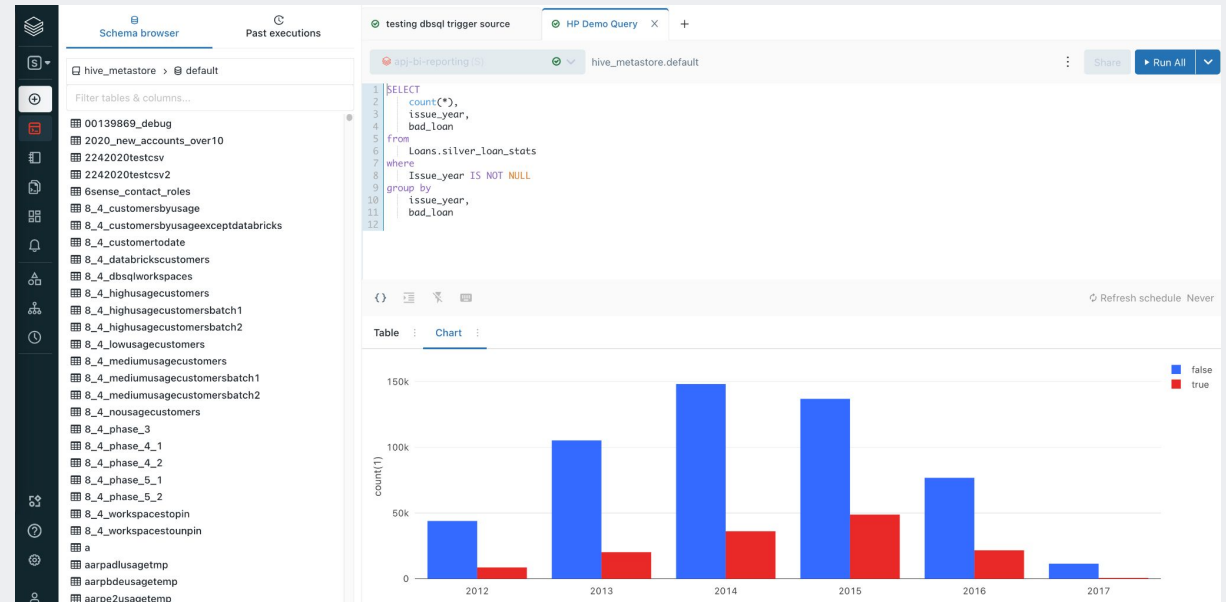
# What is serverless and why?

# Databricks SQL

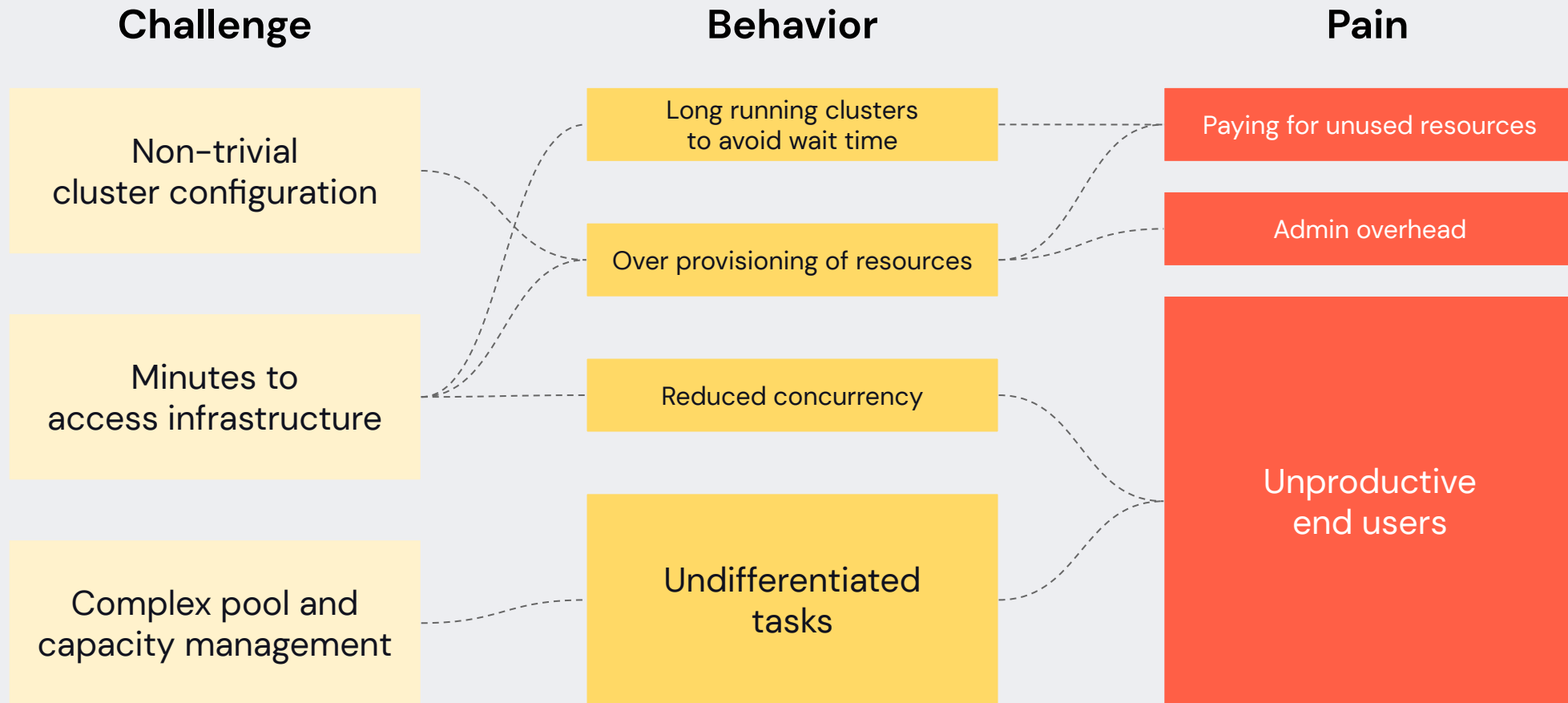## Best warehouse is Lakehouse

Databricks SQL provides an environment for:

- Running ad–hoc SQL queries
- Creating dashboards
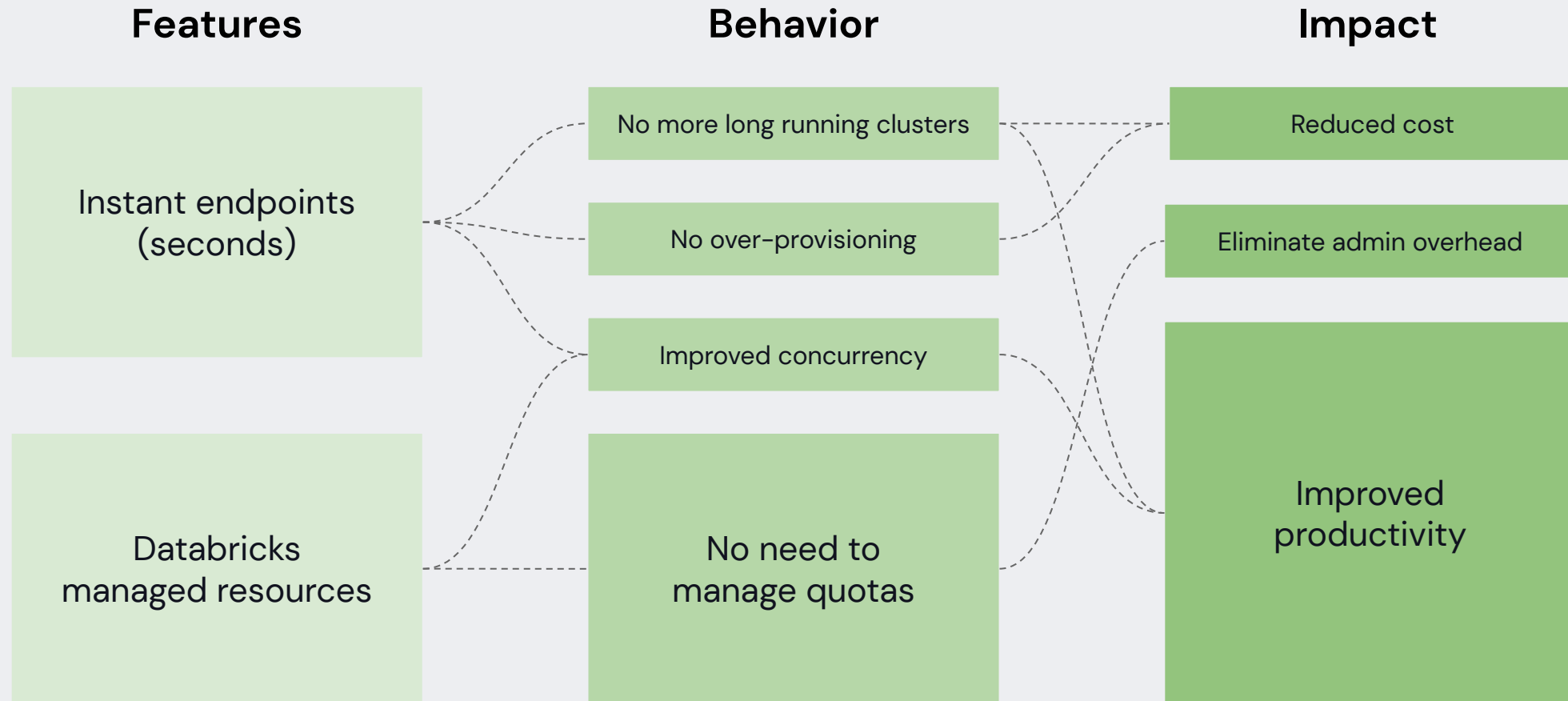- Connecting to external BI tools, e.g., Tableau and Power–BI

**Customers create and  manage the  compute cluster needed to run the queries**
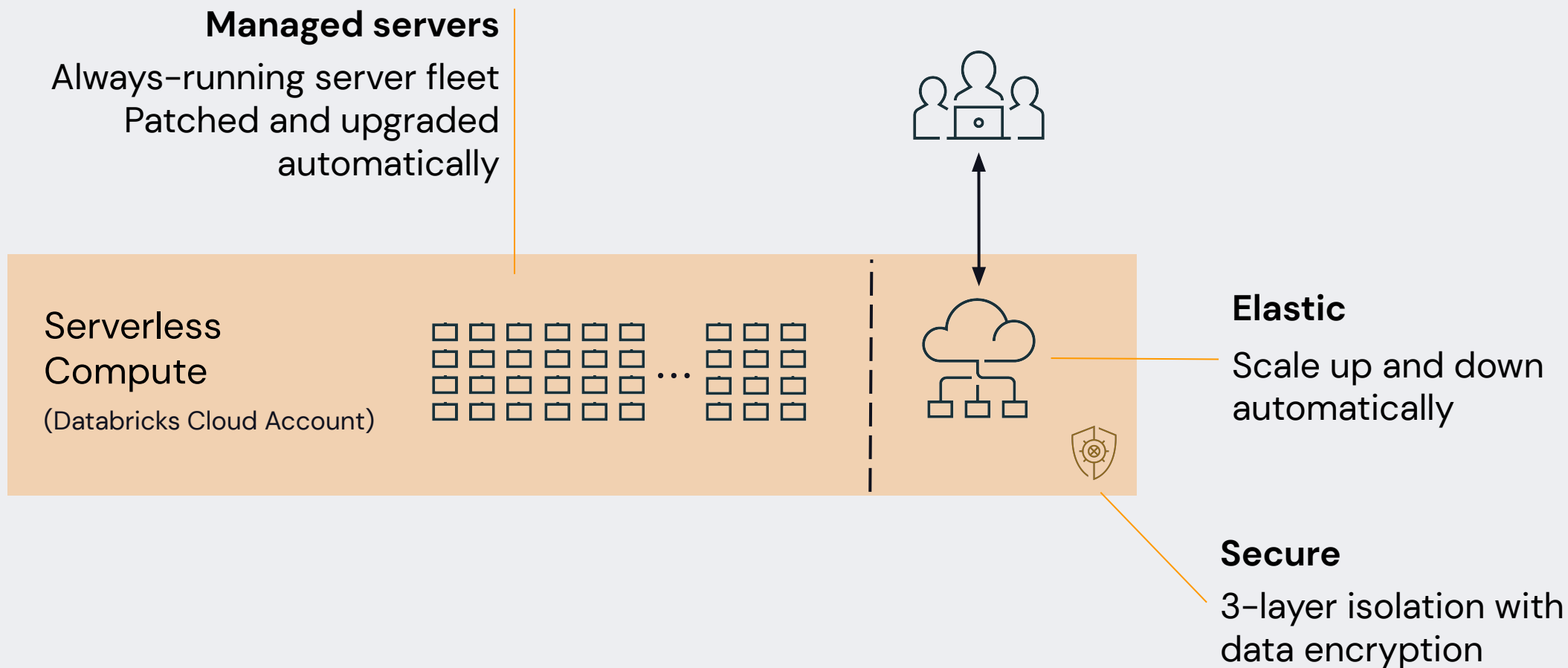
# Databricks Clusters (Classic)

| Challenge | Behavior | Pain |
|-----------|----------|------|
| Non-trivial cluster configuration | Long running clusters to avoid wait time | Paying for unused resources |
| Minutes to access infrastructure | Over provisioning of resources | Admin overhead |
| Complex pool and capacity management | Reduced concurrency | Unproductive end users |
| | Undifferentiated tasks | |

# Benefits: Databricks SQL Serverless

**Features**

**Behavior**

**Impact**

Instant endpoints (seconds)

No more long running clusters

Reduced cost

No over-provisioning

Eliminate admin overhead

Improved concurrency

Databricks managed resources

No need to manage quotas

Improved productivity

# Capabilities: Databricks **SQL Serverless**

**Managed servers**

Always-running server fleet
Patched and upgraded
automatically

Serverless
Compute

(Databricks Cloud Account)

**Elastic**

Scale up and down
automatically

**Secure**

3-layer isolation with
data encryption

# Fast 1st query performance...
# only getting better

| Classic (non-serverless) | Serverless |
|---|---|

~5 mins

~1-2 mins

Instance pools
(always on–VMs)

~10s

~2s

Current

GA

# Cost savings: Serverless platform helps reduce TCO by ~20-40%



**Compute (AWS)** ■ **DBUs**

**$37,283**

Classic SQL:
- $21,865
- $15,418
- $0.22/DBU

Serverless SQL:
- 35% customer idle cost reduction
- $24,528
- $0.70/DBU

**DATA+AI**
**SUMMIT 2022**

# Compete: Databricks Serverless SQL Serverless



Source: 2022 Cloud Data Warehouse Benchmark Report

DATA+AI
SUMMIT 2022

# How does
# it work?

# Recap: goals of Serverless

✓

## Improve reliability

- While reducing TCO

- While reducing management burden

- While maintaining (or improving) security posture

✓

## Key mechanisms:

- Instance types

- Warm pooling

- Administrative simplicity

# Instance types: classic

- Instance types critical for reliability, cost, and performance

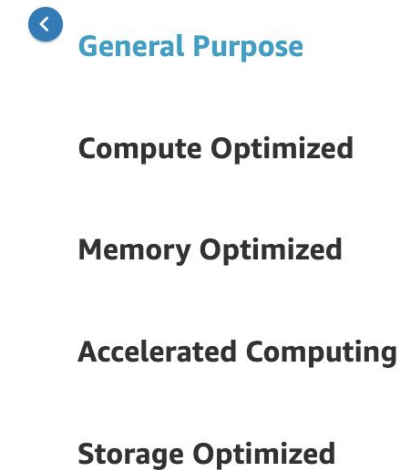| Mac | T4g | T3 | T3a | T2 | M6g | M6i | M6a | M5 | M5a | M5n | M5zn | M4 | A1 |
|-----|-----|----|----|----|-----|-----|-----|----|-----|-----|------|----|----|

**General Purpose**

**Compute Optimized**

**Memory Optimized**

**Accelerated Computing**

**Storage Optimized**

# Instance types: classic

- Instance types critical for reliability, cost, and performance

| Mac | T4g | T3 | T3a | T2 | M6g | M6i | M6a | M5 | M5a | M5n | M5zn | M4 | A1 |
|-----|-----|----|-----|----|----|----|----|----|----|----|------|----|----|

- Picking the right one(s) for a workload is hard
- Keeping up to date with new instances is hard
- Databricks can help, but...
  - Capacity
  - Reservations

**◀ General Purpose**

**Compute Optimized**

**Memory Optimized**

**Accelerated Computing**

**Storage Optimized**

# Instance types: serverless

– For SQL, we offer sizes (Small, Medium, Large, etc.)
– Leverage benchmarks to identify best cost/perf instances
– Manage capacity & reservations, enabling us to
  – Adopt new instances quickly as they become available
  – Leverage heterogeneous pool for availability
– Work with cloud vendors to understand regional/AZ capacity

# Warm pooling: classic

- Databricks offers a pooling abstraction
  - Faster startup
  - More predictable
- But: you pay for this
- Management is complicated

**Create Pool**  Cancel  Create

**Name**

**Min Idle** ?

0

**Max Capacity** ?

Optional

**Idle Instance Auto Termination** ?

Terminate instances above minimum after | 60 | minutes of idle time.

**Autopilot Option** ?

☐ Enable autoscaling local storage

**Instance Type** ?

i3.xlarge          30.5 GB Memory, 4 Cores ⌄

# Warm pooling: serverless

- Pool by default
- Amortized cost across customers
- Our responsibility to optimize pool
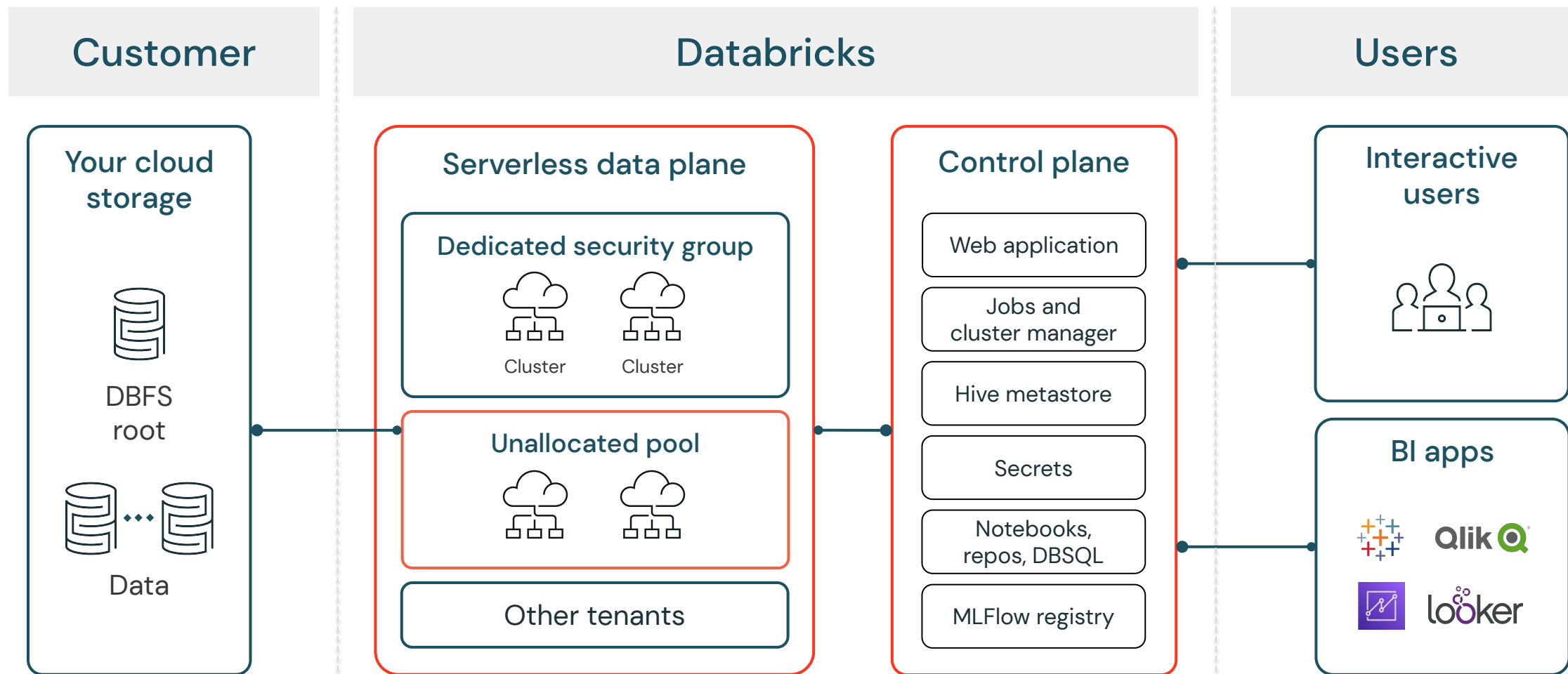
# Administrative: Classic vs. Serverless

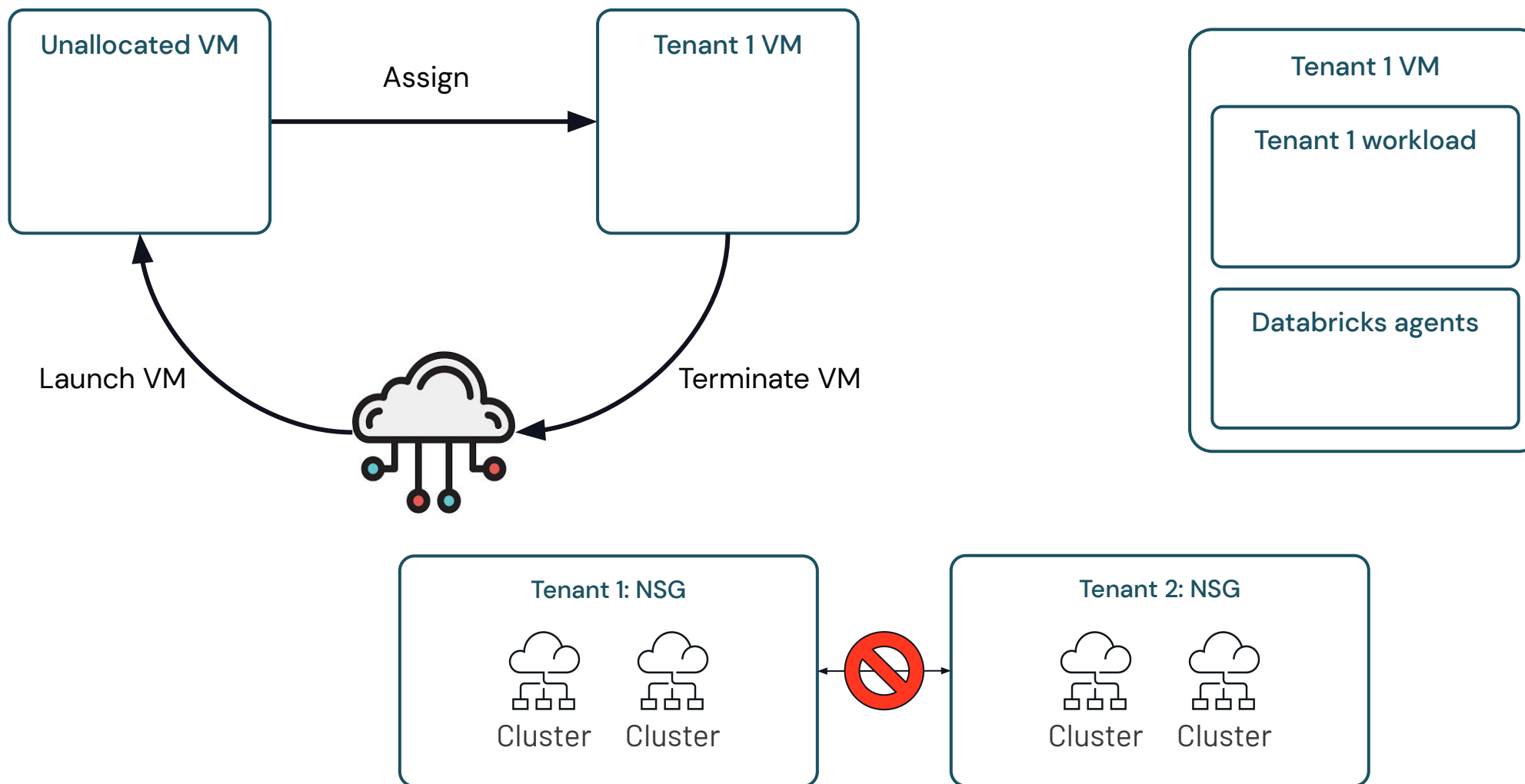- VPC/VNet CIDRs

- Account/subscription sharding

# Security

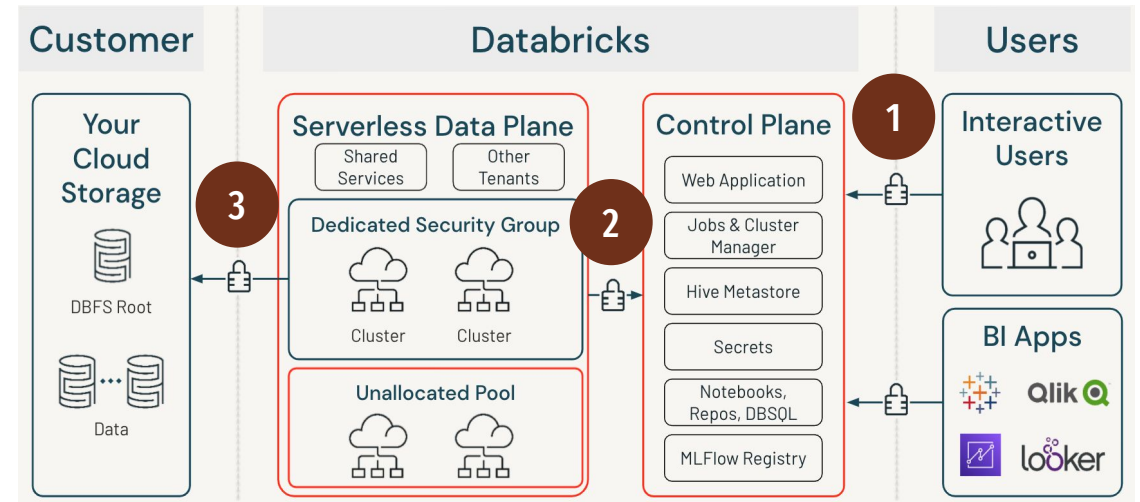# Classic architecture

# Serverless architecture

| Customer | Databricks | Users |

**Your cloud storage**

DBFS root

Data

## Serverless data plane

### Dedicated security group

Cluster          Cluster

### Unallocated pool

### Other tenants

## Control plane

Web application

Jobs and cluster manager

Hive metastore

Secrets

Notebooks, repos, DBSQL

MLFlow registry

**Interactive users**

**BI apps**

# Serverless VM-level security

Unallocated VM → **Assign** → Tenant 1 VM

Launch VM

Terminate VM

**Tenant 1 VM**
- Tenant 1 workload
- Databricks agents

**Tenant 1: NSG**
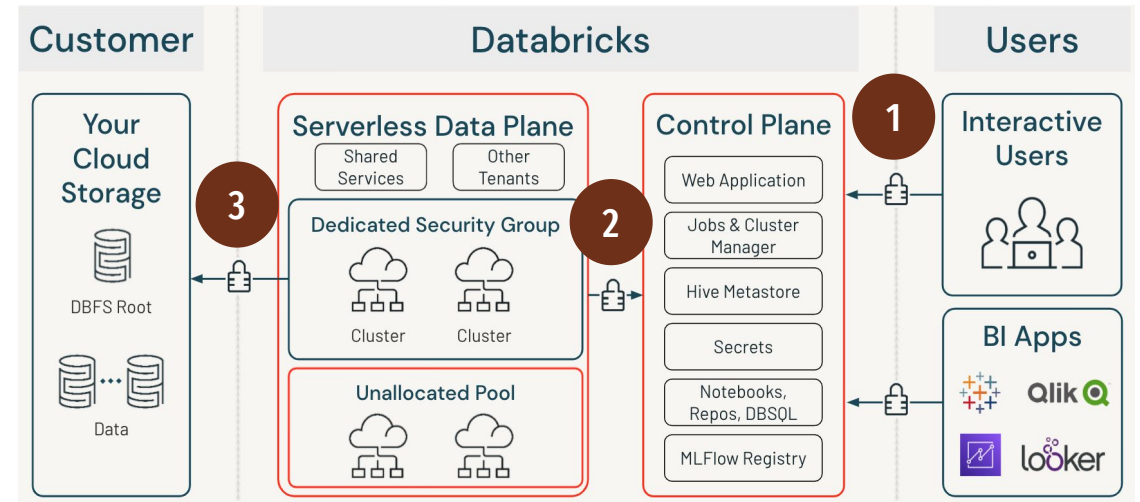- Cluster
- Cluster

**Tenant 2: NSG**
- Cluster
- Cluster

# Network access

1. User to Control Plane Options
   - Open to all (Typical)
   - **IP Access Lists**
   - **Private Link** (Preview)

# Network access

1. User to Control Plane Options
   - Open to all (Typical)
   - **IP Access Lists**
   - **Private Link** (Preview)

2. Data Plane to Control Plane managed by Databricks
   - mTLS 1.2+
   - **All workers are private**

# Network access

1. User to Control Plane Options
   - Open to all (Typical)
   - **IP Access Lists**
   - **Private Link** (Preview)

2. Data Plane to Control Plane managed by Databricks
   - mTLS 1.2+
   - **All workers are private**

3. Data Plane to Customer Storage
   - Assume Role via Public IP
   - **Private network connectivity to blob storage**

# Three layers of isolation controls

**1. Container Isolation** **2. VM Isolation** **3. Network Isolation**

- Hardened container images per industry best practice

- Disable privileged access in the container

- No VMs reuse

- No privileges within broader environment

- Federated access via temporary security tokens

- Only intra–cluster traffic allowed

Tested by internal teams and external vendors

# Serverless Resource Life-cycle

# Features and limitations

- Internal Hive Metastore is **supported**

- External Hive Metastore access via public IP is **supported**

- Unity Catalog is **supported**

- Private Link (DP -> S3/ADLS)

- Python/Scala are UDF not supported

- 5min auto-stop and 1min via API **in Q2**

# Data location and encryption

- Data <u>stays</u> in customers cloud account (S3 buckets)

- Databricks managed Data Plane is in the same regions customers data plane—**no egress cost**

# Serverless Access to Big Healthcare Data

KYTHERA™
DECIPHERING HEALTHCARE

# KYTHERA LABS & DATABRICKS

## Provided via a Big Data Platform

‹ Databricks' First healthcare OEM

‹ Ability to spin up and manage branded instances for prospects and clients (Wayfinder)

‹ Wayfinder gets clients access to refined data assets within hours - not months

‹ Increases return on data investments, 8x faster results, and richer, more granular insights

## We Provide Remastered Healthcare Data and Analytics

‹ 330+ million patients

‹ 12.5 billion healthcare claims

‹ 12.9 billion Rx claims

‹ 27 billion claim lines

‹ BYO Data

## With Rich Analytics and Expertise

‹ Unified and simplified data model

‹ Corrected, standardized and flattened data models (Remastered)

‹ Foundational patient pathway analytics

‹ Machine learning models

‹ BI Visualizations

‹ Useful interfaces (Datavant, Salesforce, Mirador, etc)

KYTHERA

# DISCOVER TARGET POPULATIONS

Healthcare and Life Sciences companies often need to identify and validate patient populations and markets to understand how patients behave over time - where they receive care, who provides care, and how care is reimbursed - to make informed business decisions.

Patient Cohort dashboards empower business users to access and analyze data without wading through billions of rows of data, storing idle data, or waiting for a cluster to start.

KYTHERA

# *SERVERLESS COHORT BUILDER*

## *Faster access to solve the toughest data challenges*

KYTHERA

# BETTER INSIGHTS. FASTER.

Biggest Pain Points Addressed by Serverless
- ‹ Perfect for casual access by citizen analyst
  - ○ Infrequent access to data, but need for quick response time
  - ○ Serverless cuts out the cluster startup time
- ‹ Manage costs
  - ○ Sure, with lots of money we could just have a big cluster waiting for that infrequent access
  - ○ Serverless eliminates the cost of standby compute by leveraging shared resources across the customer base
- ‹ Accelerate the time to value and illuminate critical insights

Contact:
Matt@kytheralabs.com
Kytheralabs.com

KYTHERA

# THANK YOU

*matt@kytheralabs.com*

*@kytheralabs*

*linkedin.com/company/kytheralabs*

KYTHERA

# Roadmap

# Roadmap: serverless platform



DBSQL Serverless Public preview

DBSQL Serverless GA

7/1

NOW

10/1

1/1

4/1

6/1

Serverless: Notebook and Jobs (previews)

Serverless: Notebook and Jobs (GA)

CY 2023

# Try it now!

1. Databricks SQL customers on AWS:

   a) Enable Serverless from account console
   b) Upgrade all endpoints from classic
      to Serverless

2. Databricks SQL customers on Azure:

   a) Submit your interest:
      bit.ly/DBSQLServerless-Azure

3. Learn more, talk to your
Databricks representative

**DATA+AI**
SUMMIT 2022

Thank you