As the title suggests, this is a demo heavy session.
The slides are more for the talktrack, which will be pretty minimal.
Will potentially have the intro slides and then will jump into demos.

DATA⁺AI SUMMIT
BY databricks

# DATABRICKS WORKFLOWS: PRACTICAL HOW-TOS AND DEMOS

**Author Name**
**Date**
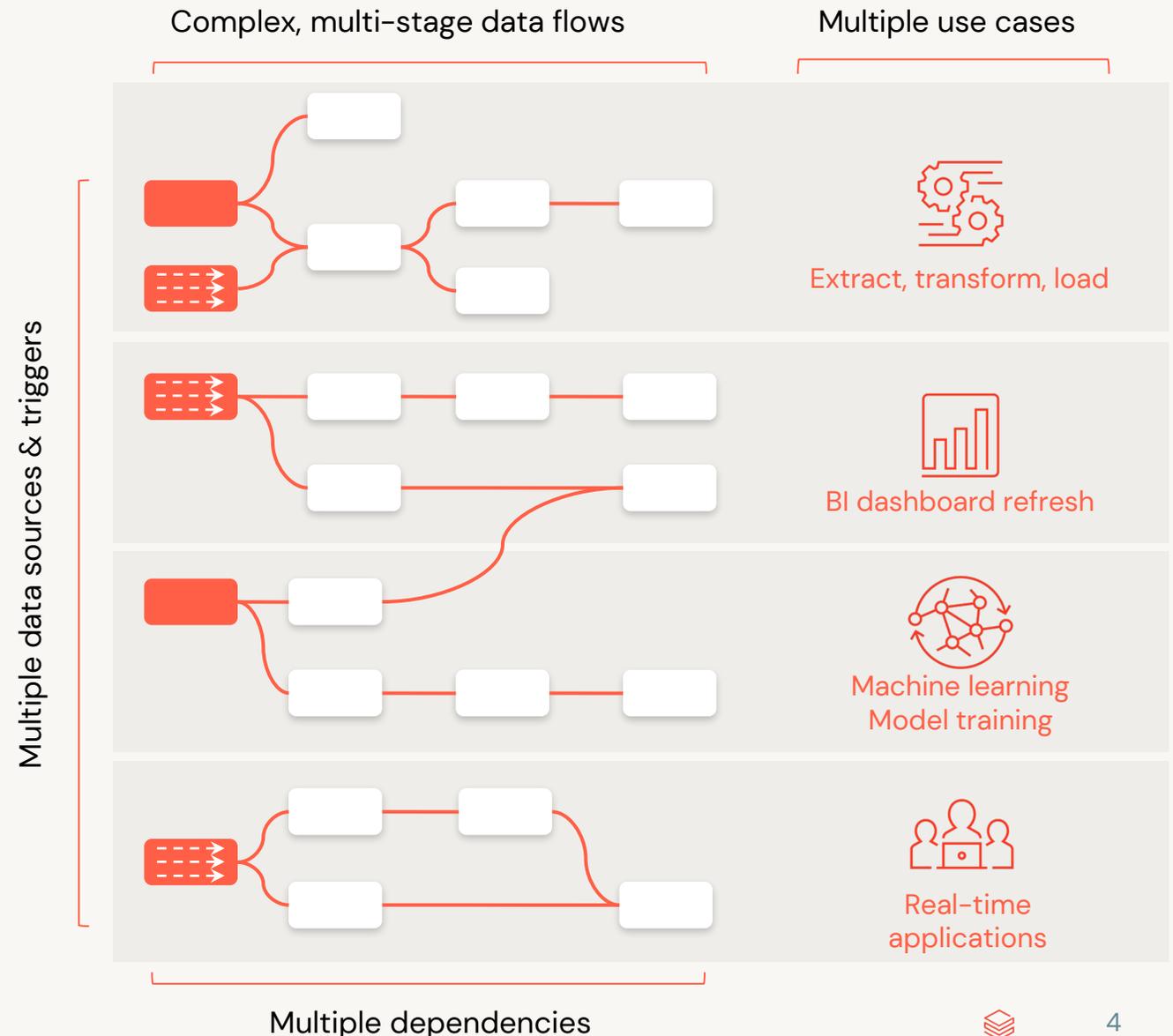
# Product safe harbor statement

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

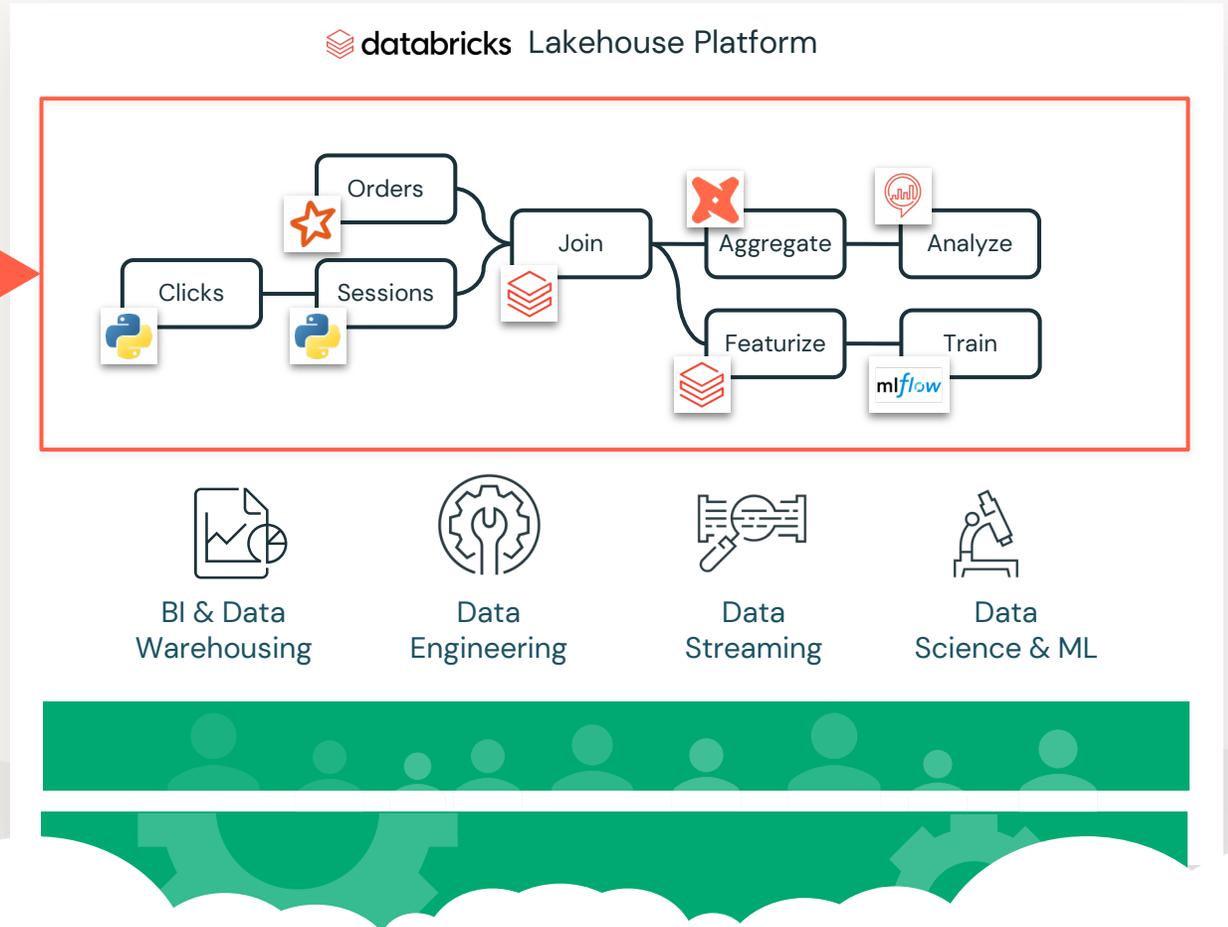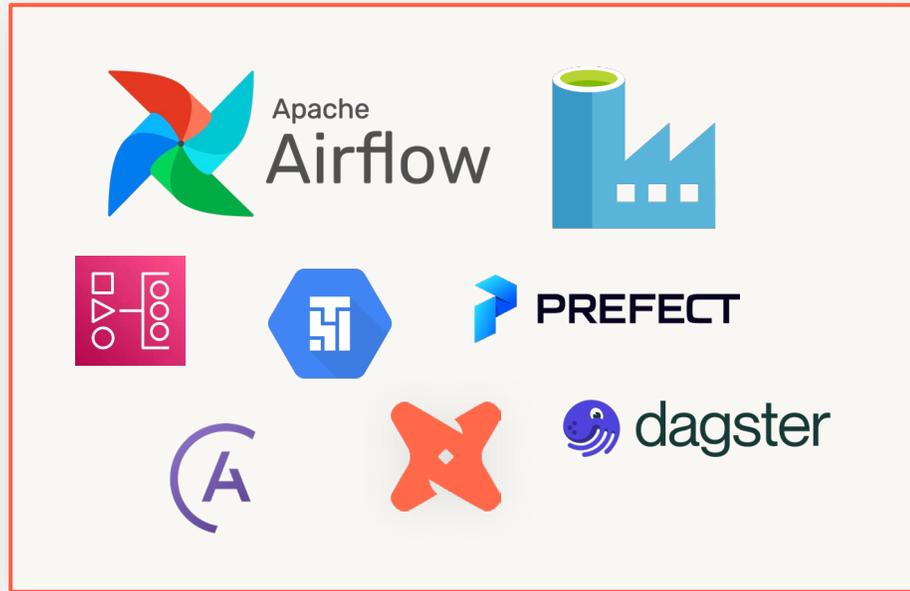# Modern data engineering requires modern data orchestration

**Orchestrating processes across all data, analytics and AI use cases is business critical**

"Data pipelines are growing in size, volume, and complexity, with multistage processing and dependencies between various data assets."*

*Gartner Data Engineering Essentials, Patterns and Best Practices, September 2022*

Complex, multi-stage data flows

Multiple use cases

Multiple data sources & triggers

Multiple dependencies

Extract, transform, load

BI dashboard refresh

Machine learning Model training

Real-time applications

# There are many ways to orchestrate your Lakehouse workloads

# External orchestrators create challenges

| Hard to use for many practitioners | Difficult to understand root cause when issues occur | Complex architecture to manage and maintain |
|---|---|---|
| Data teams are less productive | Bad data lowers value of downstream applications | Higher cost of ownership and lower reliability |

Apache Airflow

**These tools are not unified with your Lakehouse**
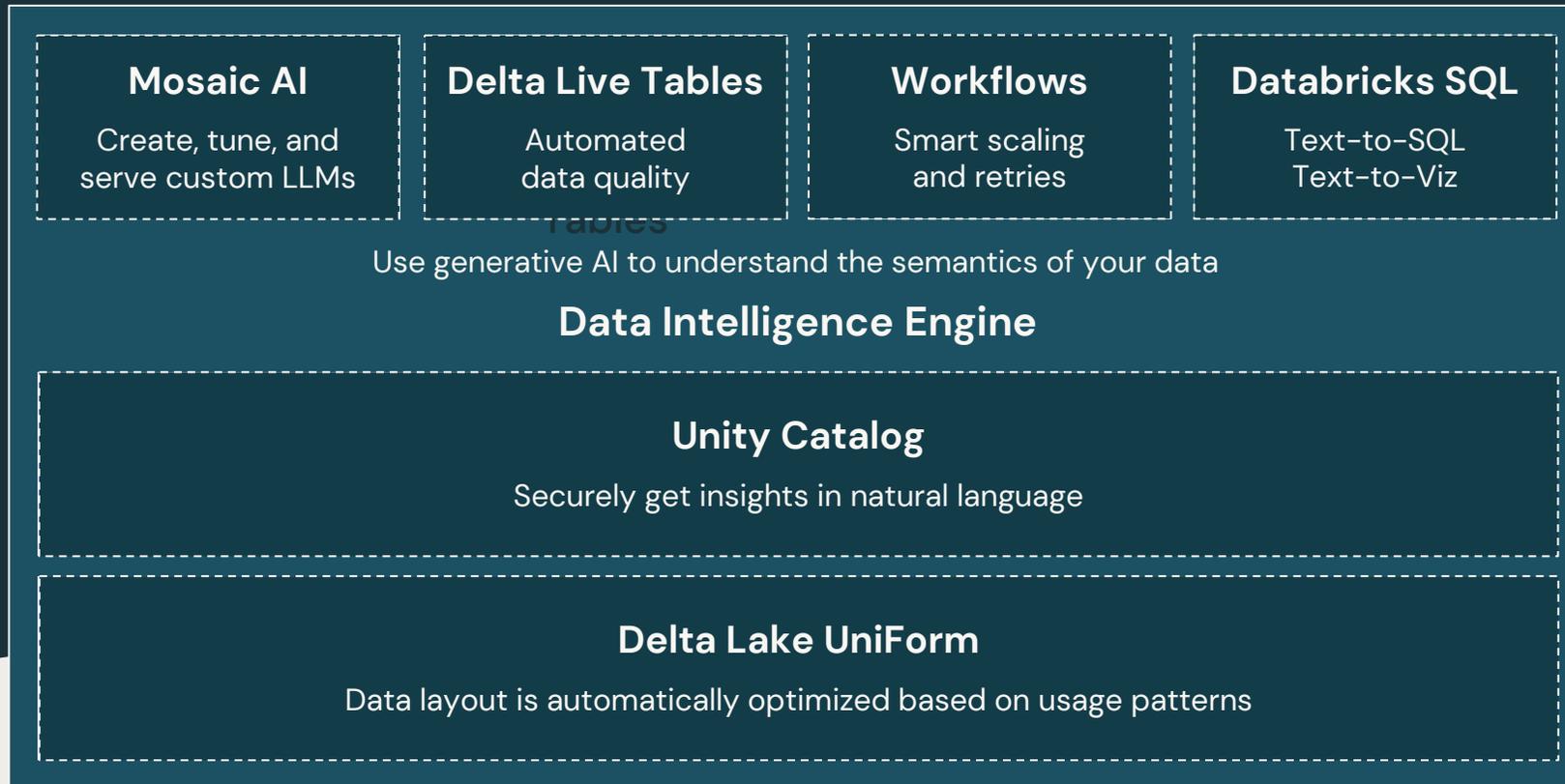
dagster

PREFECT

# Databricks Workflows

## Unified orchestration for data, analytics, and AI on the Lakehouse Platform

- Simple authoring
- Actionable insights
- Proven reliability

# Databricks Data Intelligence Platform

**Mosaic AI**
Create, tune, and serve custom LLMs

**Delta Live Tables**
Automated data quality

**Workflows**
Smart scaling and retries

**Databricks SQL**
Text-to-SQL
Text-to-Viz

Use generative AI to understand the semantics of your data

**Data Intelligence Engine**

**Unity Catalog**
Securely get insights in natural language

**Delta Lake UniForm**
Data layout is automatically optimized based on usage patterns

**Open Data Lake**

All Raw Data
(Logs, Texts, Audio, Video, Images)

# Top 3 reasons why customers love Databricks Workflows

## Simple authoring
### for all data practitioners

Any data practitioner accelerate development by easily orchestrating workflows from inside their Databricks workspace in just a few clicks. Advanced users can use their favorite IDE's with full support for CI/CD.

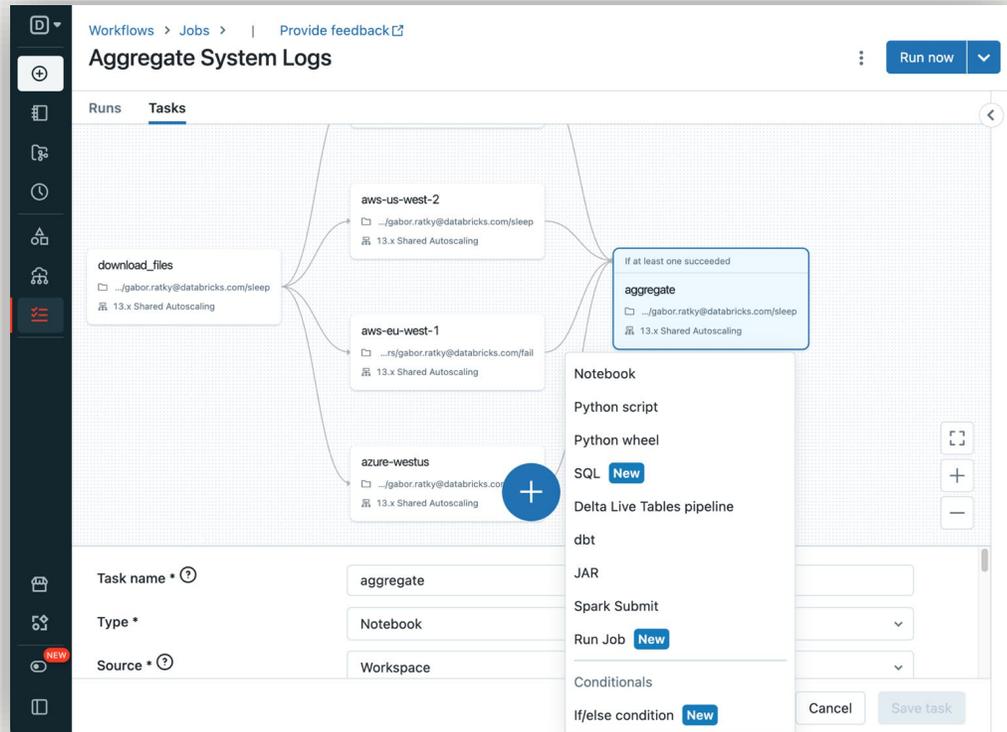## Actionable insights
### from real–time monitoring

Full visibility into every task in every workflow. See the health of all your production workloads in real–time with detailed metrics and analytics to identify, troubleshoot, and fix issues fast.

## Proven reliability
### for production workloads

A fully managed orchestration service with serverless data processing and a history of 99.95% uptime. Workflows is trusted by thousands of Databricks customers running millions of production workloads.

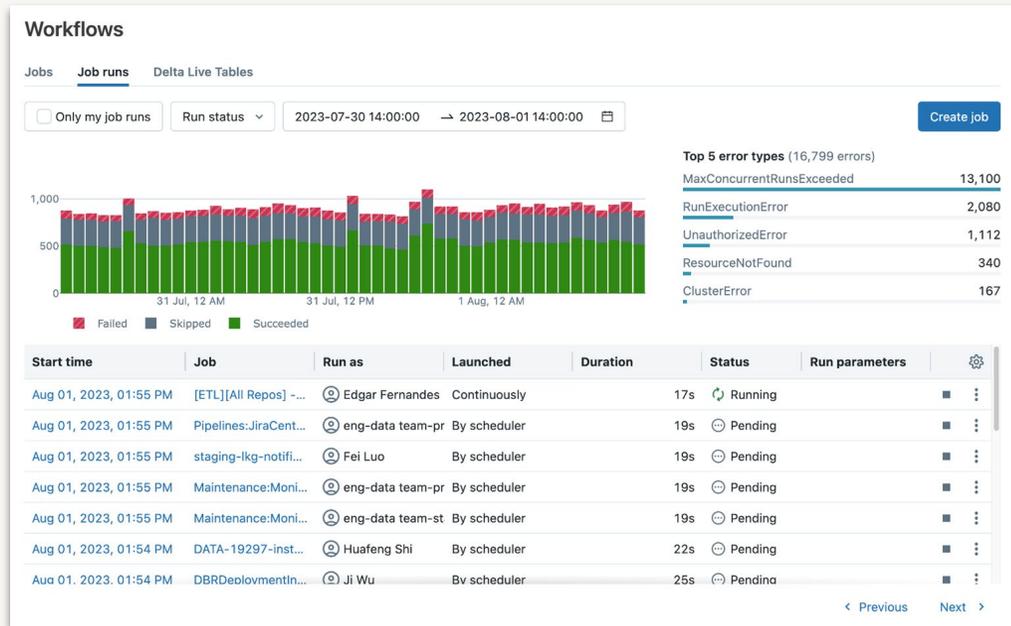# Simple authoring for all data practitioners



Build sophisticated workflows inside your Databricks workspace with a few clicks
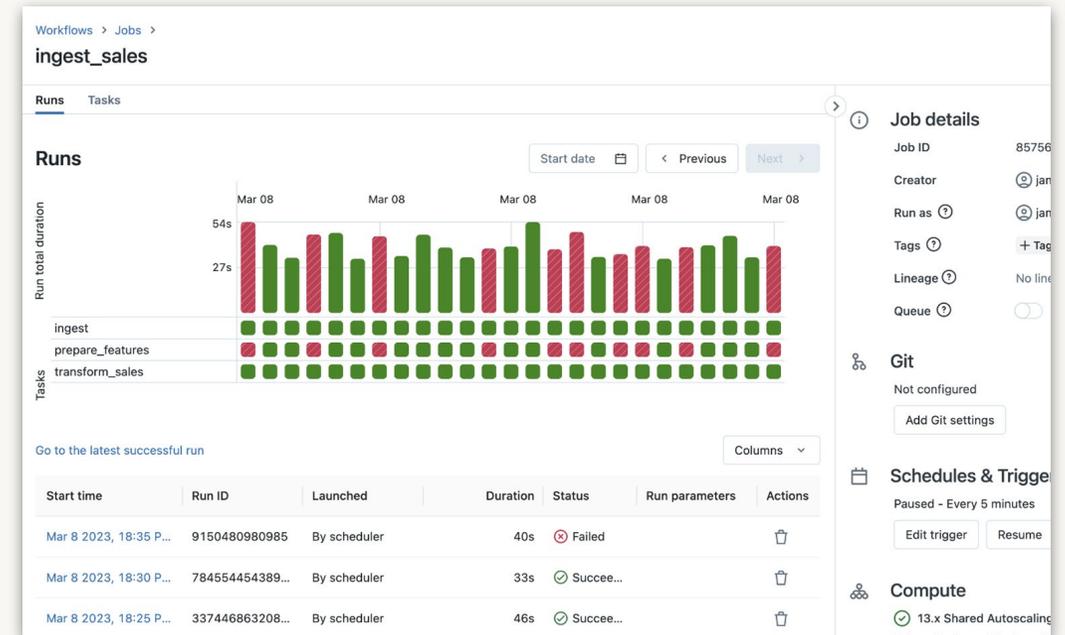
Or connect your favorite IDE to develop workflows locally and run them on Databricks

# Actionable insights from real-time monitoring



A simple and intuitive monitoring UI provides real-time metrics and detailed analytics for every workflow run

Drill down to understand which tasks are failing and why. Troubleshoot issues before your customers are impacted

# Proven reliability for production workloads

## 99.95%uptime

Trusted by thousands of customers running millions of production workloads

### Fully managed service
Reduce maintenance costs and let your teams focus on innovation instead of resource management tasks

**[In Preview]**
### Serverless data processing
Massively scalable compute resources for task execution remove even more of the maintenance burden from your teams and further reduces costs

# Building Blocks of Databricks Workflows

A unit of orchestration in Databricks Workflows is called a **Job.**

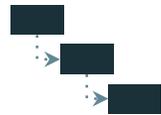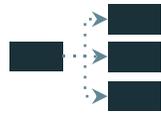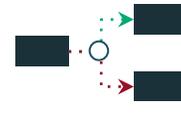| **Jobs** consist of one or more **Tasks** | Databricks Notebooks | Python Scripts | Python Wheels | SQL Files/Queries | DBSQL Dashboards | Delta Live Tables Pipeline | dbt | Java JAR file | Spark Submit |
|---|---|---|---|---|---|---|---|---|---|

| **Control flows** can be established between **Tasks.** | Sequential | Parallel | Conditionals (Run If) | Jobs-as-a-Task (Modular) | For-Each Loop |
|---|---|---|---|---|---|

| **Jobs** supports different **Triggers** | Manual Trigger | Scheduled (Cron) | API Trigger | File Arrival Triggers | Table Triggers *Coming Soon* | Continuous (Streaming) |
|---|---|---|---|---|---|---|

# More on Operations and Cost Efficiencies

## Resource Isolation

**Dedicated, tailored** Job clusters enables each task to run independently **without resource sharing bottlenecks!**

## Job Cluster Re-Use

**Cluster re-use** enables users to run tasks in a Databricks Job on the same cluster for more **efficient cluster utilization and decreased job latency!**

## Late Running Jobs

Tasks can **trigger emails** to stakeholders enabling you to **proactively monitor and take action against late or long running jobs!**

## Repair and Re-Run

You can configure **when and how many times failed runs are retried!**

# Serverless Compute

## SIMPLE and FAST

No knobs
Fast startup
For any practitioner

## EFFICIENT

Fully managed and versionless
Paying only what you use
Strong cost governance

## RELIABLE

Secure by default
Stable with smart fail-overs

**GA**

| DB SQL | Workflows | Notebooks | Delta Live Tables |

## Serverless Compute

Hands-off auto optimized compute managed by Databricks

Storage

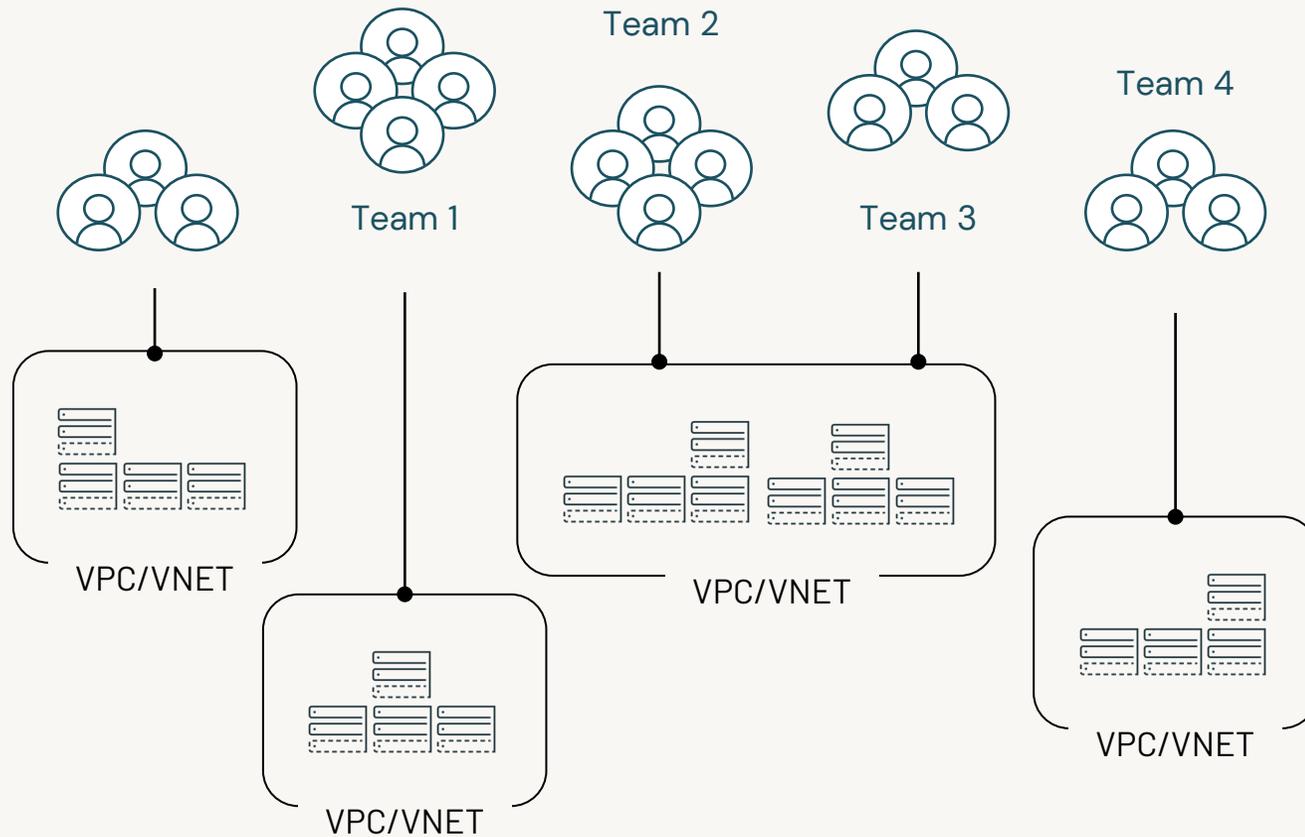# STOP spending time on…

### Setting up networks

Create and configure VNets
Set up gateways and firewall rules
Setup and manage private endpoints
X-tenant identities
IP address / subnet management

### Security and Compliance

Vulnerability management
Encryption and key management
Intrusion detection and monitoring Data
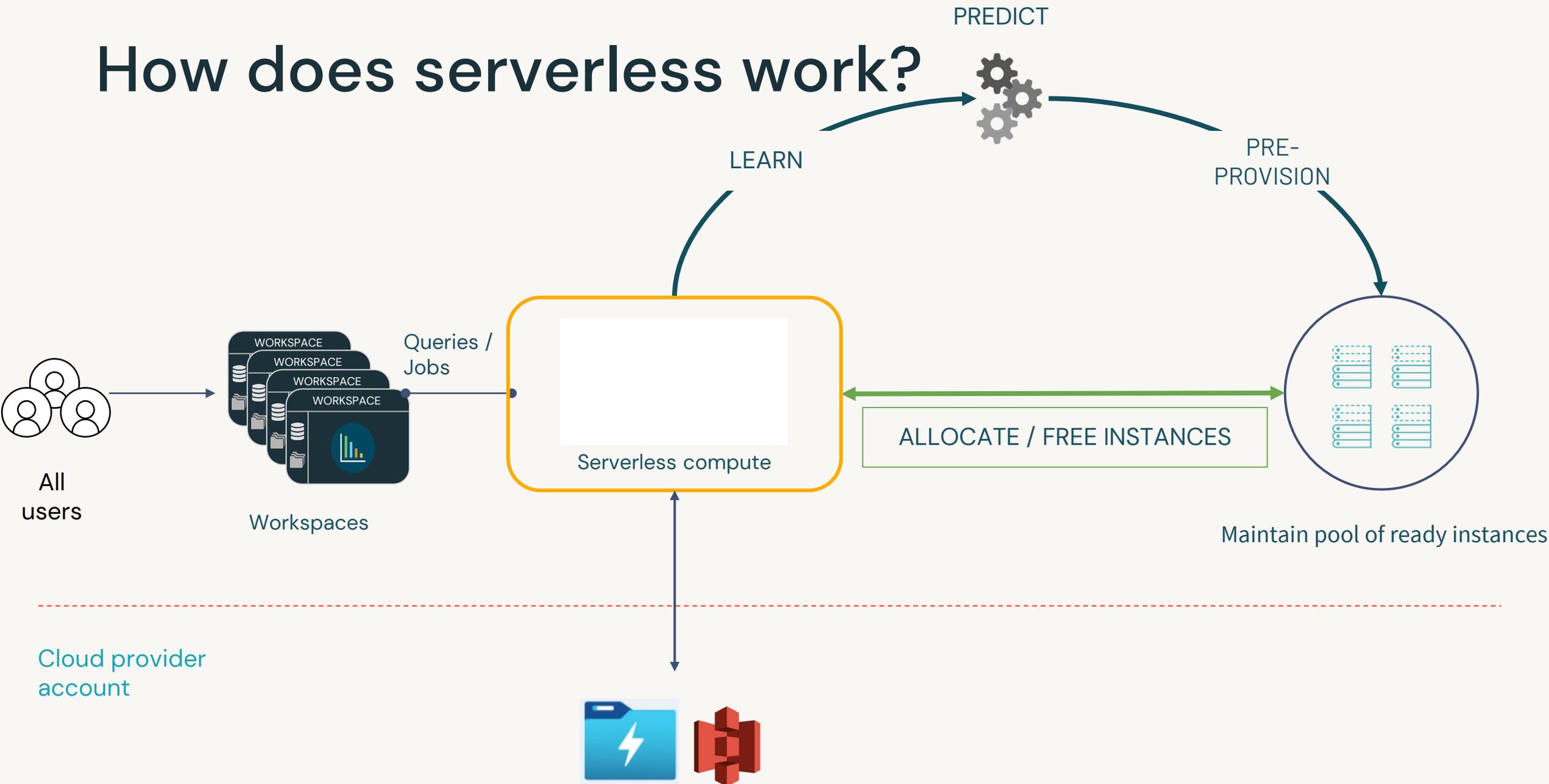exfiltration protection

### Managing efficiency

Capacity projections and reservations
Right sizing instances for workloads
Maintaining high utilization
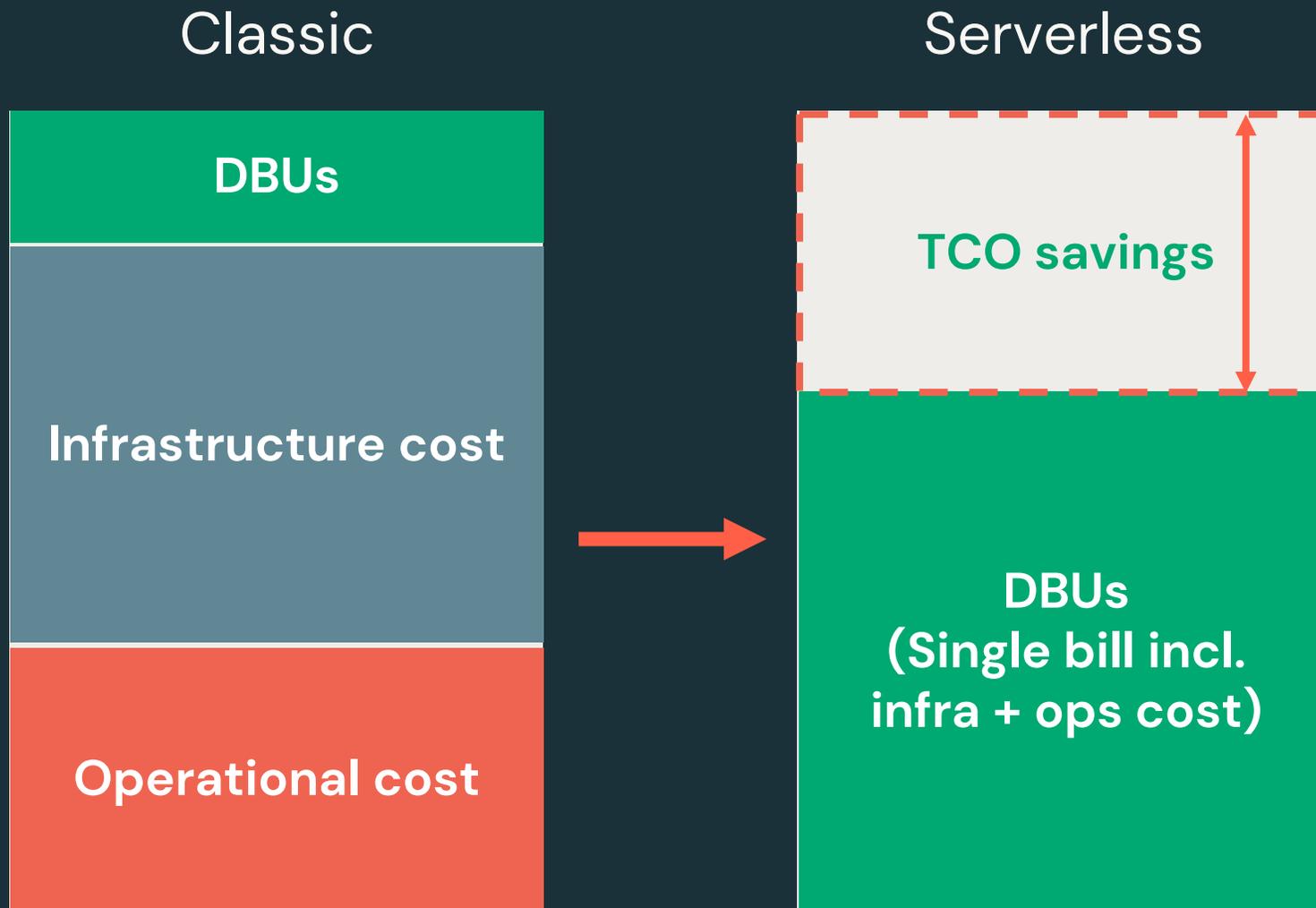Managing instance pools
Vacuum / compaction of Delta tables

Team 1
Team 2
Team 3
Team 4

VPC/VNET
VPC/VNET
VPC/VNET
VPC/VNET

**Structured, Semi-structured and Unstructured Data**

# How does serverless work?

PREDICT

LEARN

PRE-PROVISION

All users

Workspaces

Queries / Jobs

Serverless compute

ALLOCATE / FREE INSTANCES

Maintain pool of ready instances

Cloud provider account

18

# Serverless reduces TCO

**Classic**

**Serverless**

DBUs

Infrastructure cost

Operational cost

TCO savings

DBUs
(Single bill incl.
infra + ops cost)

# Serverless Compute for Workflows

## Hands-off, auto-optimizing compute

**Fully managed and reliable**

- <60s startup (at GA)
- Versionless with auto-update
- System tables
  for cost observability
- Budgets for cost control
  (coming)
- Automatic instance type
  failover (coming)

# Secure, multi-user, serverless Spark

## Notebooks with serverless compute



## Workflows with serverless compute



**Built on Spark Connect:** user code runs in full isolation on client, driver and executors

| REPL<br>Python code (non-Spark) | | Spark<br>Connect | | Spark Driver | | Spark Executors |
|---|---|---|---|---|---|---|
| Client Isolation | | | | Driver Isolation | | Executor Isolation |

# HowTos and Demos with emphasis on best practices