# DATA⁺AI SUMMIT
BY databricks

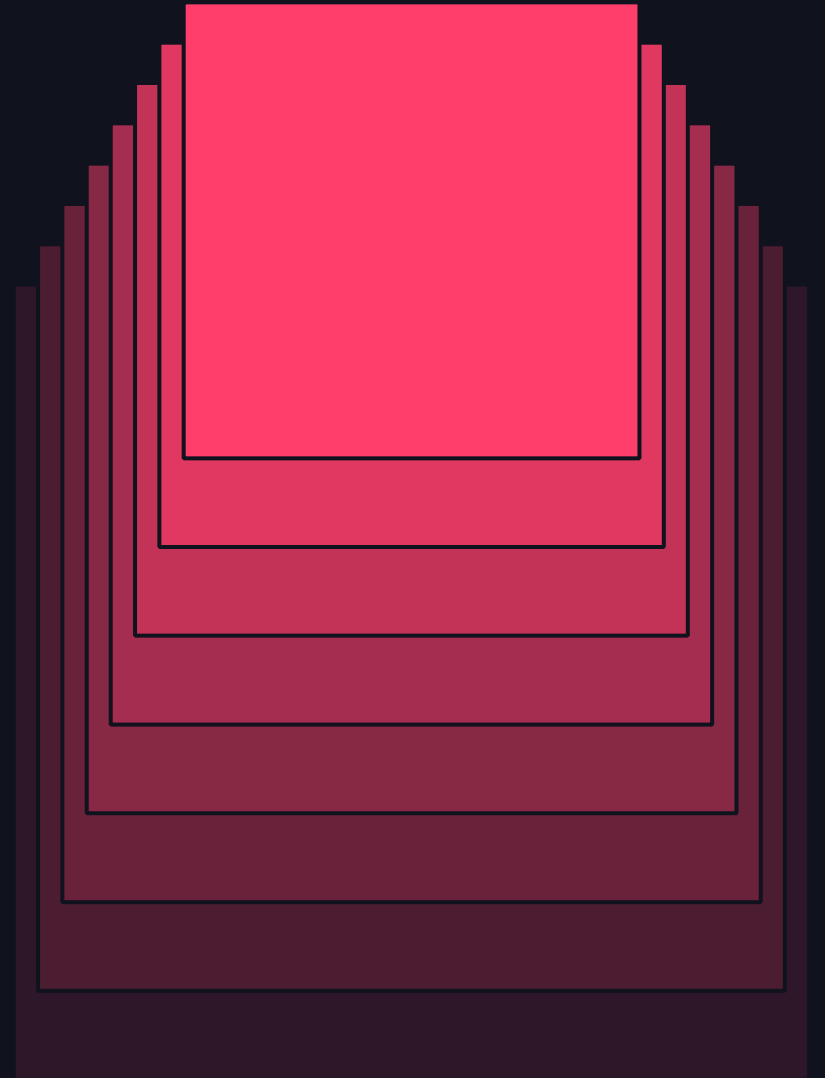# What's new in Databricks Workflows

R.R. Fäustlin, Product Management, Databricks
June 12, 2024

DATA⁺AI SUMMIT

# Product safe harbor statement

# Agenda

- The Databricks Workflows Story
- Recent innovations
- Looking ahead
- Demo

- **2015 Cron-based jobs**
- **2016 Notebook workflows**
- **2020 Jobs with multiple tasks**
  - Reliability
  - Monitoring
- **2022 The best lakehouse orchestrator**
  - Integration with the lakehouse
  - Streaming
  - Cluster reuse
- **2023 Serverless, performance and ease of use**
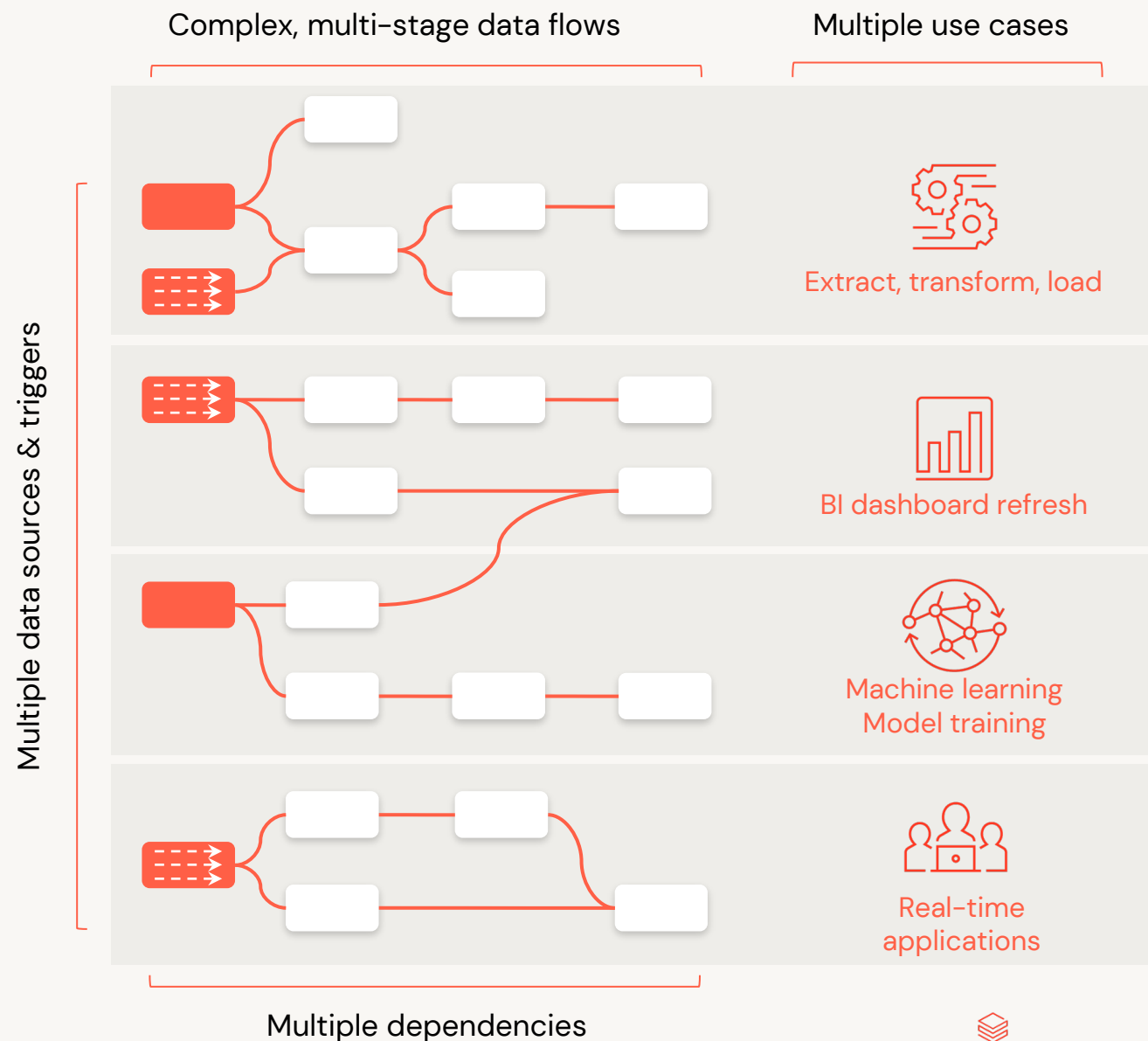- **2024 AI powered ETL**

# Modern data engineering requires modern data orchestration

# 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

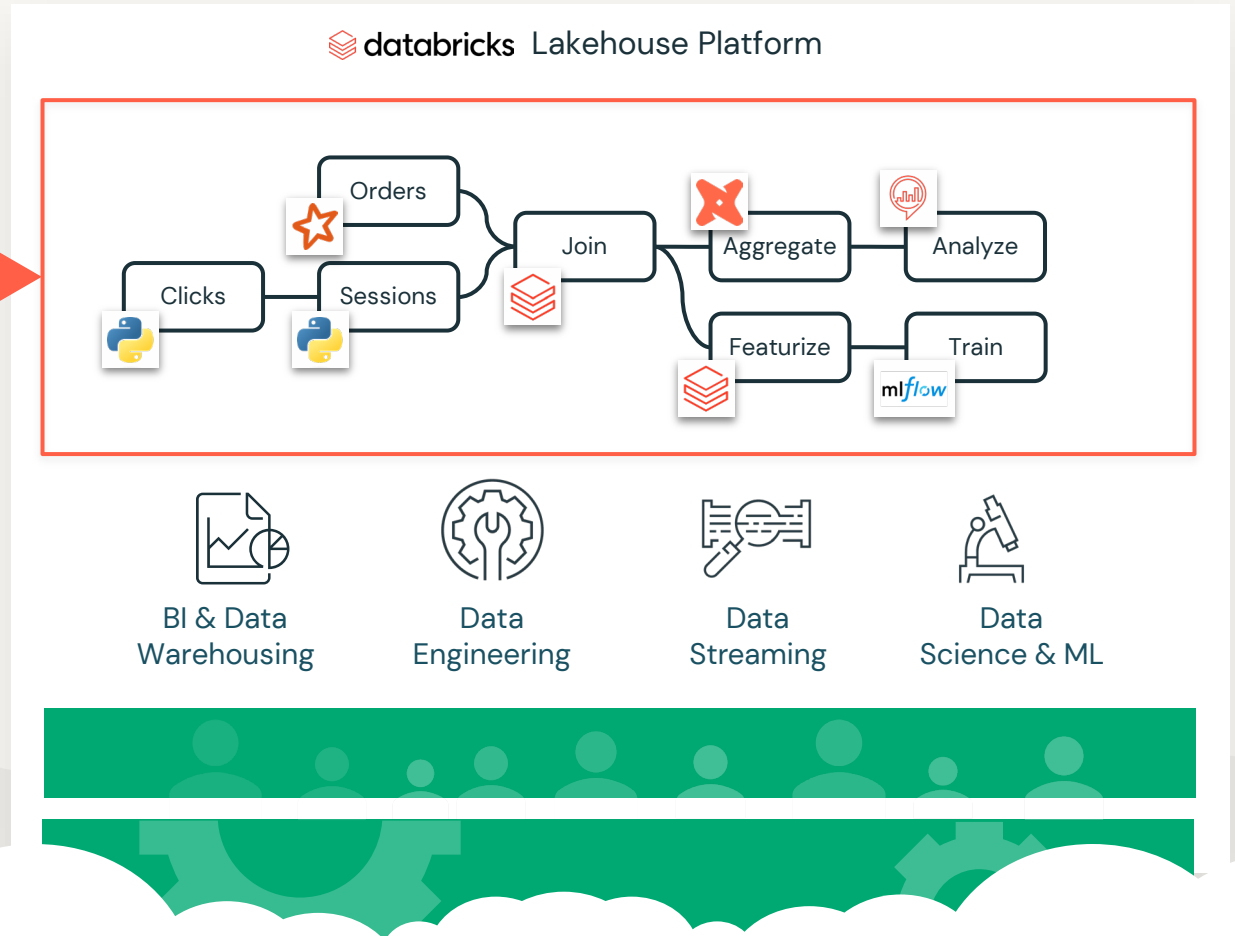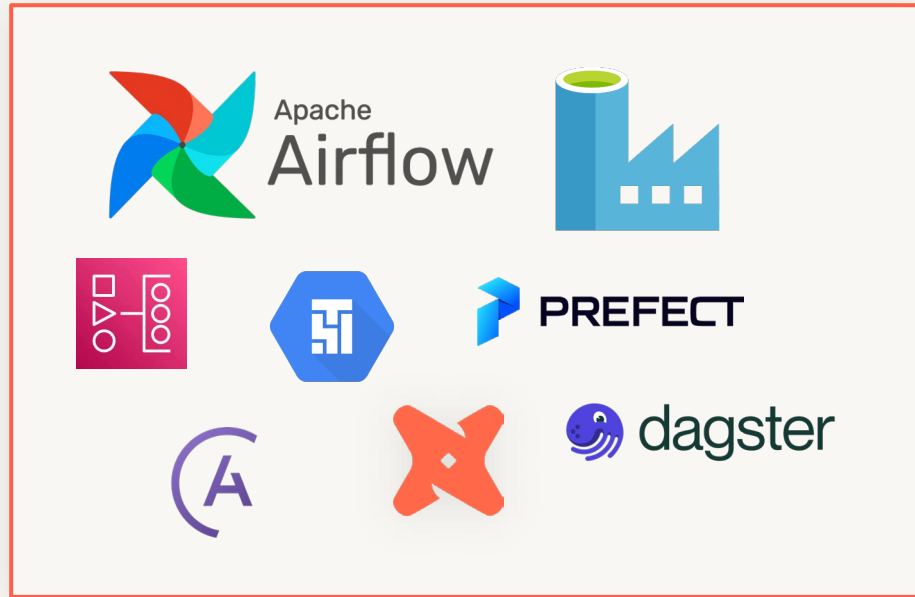# But organizations struggle with so many tools

**>65%**

of organizations are using 10+ data engineering and intelligence tools

Source: IDC DataOps Survey, 2020

# Many ways to orchestrate your Lakehouse

# External orchestrators create challenges

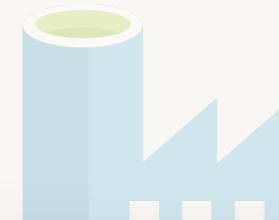| Hard to use for many practitioners | Difficult to understand root cause when issues occur | Complex architecture to manage and maintain |
|---|---|---|
| ↓ | ↓ | ↓ |
| Data teams are less productive | Bad data lowers value of downstream applications | Higher cost of ownership and lower reliability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

Apache Airflow

dagster

PREFECT

## These tools are not unified with your Lakehouse

# Databricks Workflows

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

- **Simple authoring**
- **Actionable insights**
- **Proven reliability**

# Top 3 reasons why customers love Databricks Workflows

## Simple authoring
**for all data practitioners**

**Any data practitioner can** accelerate their development by easily orchestrating Workflows from inside their Databricks workspace in just a few clicks. Advanced users can use their favorite IDEs 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** service with serverless data processing and years of **99.95% uptime**. Workflows is trusted by thousands of Databricks customers running millions of production workloads.

**>10k** customers | **>25 million** VMs/day | **>99.95%** uptime

Disney streaming services    Adobe    AT&T    ABN·AMRO    John Deere    COMCAST

FedEx    zalando    Walmart    Nielsen    HSBC

Wood Mackenzie
A Verisk Business

**Improved collaboration
80-90% faster processing**

Ahold Delhaize

**ADF → Workflows
4.5x faster deployment
50% cost reduction**

yipit DATA

**Airflow → Workflows
60% cost reduction
90% faster processing**

# >70 new features 🚀 shipped 🚀

# This is the highlight reel.

# Data triggers: Run only when you need to

✅ Trigger based on table change

**SOON** Trigger when new files arrive
Now unlimited files count and
increasing number of triggers

✅ Trigger another job
"job as a task"

---

**Schedules & Triggers** ✕

**Trigger Status**
⦿ Active
○ Paused

**Trigger type**
Table update ▾

Table update triggers monitor tables for data changes (e.g. update, merge and delete). These tables can be managed or external tables in Unity Catalog.

**Tables** ⓘ

| mycatalog.myschema.mytable1 | 🗑 |
| mycatalog.myschema.mytable2 | 🗑 |
| mycatalog.myschema.mytable3 | 🗑 |

**+ Add table**

**Trigger when** ⓘ
⦿ All tables are updated
○ Any table is updated

**Advanced** ⌄

Test connection    Cancel    Save

# Advanced SQL orchestration

✅ Referencing of SQL query results in other tasks, e.g. for conditional execution

✅ Multi-SQL statement support

→ Full support for control flow, e.g. conditionals, for-each

# Faster and easier

✅ Now up to 1k tasks per job

✅ Improved cluster defaults to set you up for success

✅ Easily exchange messages across tasks, now with a simplified UI and auto-complete

# Find your assets quickly

✅ Easily filter to the job or run that you care about, e.g. with job favourites ⭐

✅ Descriptions for your jobs and tasks

# ✅ AI assisted debugging integrated

# Only run what you need

✅ Run only the tasks you need to get back on track, now also single and successful tasks

# Easily optimize price/performance

✅ Timeline view across tasks and queries

✅ Query profile integration

✅ Track streaming lag and alert on deviations

✅ Alert when jobs are running late

# Track cost and long–term trends

✅ System tables and templated/customizable dashboards

**SOON** In–UI budget monitoring and alerting

# PyDABs: Anything in Databricks as code

☑ Python SDK

☑ Terraform support

☑ Run jobs as service principal

→ Easily develop Workflows in your IDE as Python code

→ Compare changes

→ Collaborate with UI-only users

## Code review your Workflows

# INGESTION CONNECTORS:

Efficient data ingestion for everyone

**Simple** and low-maintenance ➡ Fewer headaches, quicker time to value, democratized data

**Unified** with the lakehouse ➡ Secure and healthy pipelines that live where you do your work

**Efficient** end-to-end ➡ Lower costs, better performance, better scalability

| salesforce | workday. | SQL | ORACLE NETSUITE | servicenow | Google Analytics | SharePoint | PostgreSQL |
|---|---|---|---|---|---|---|---|
| Private preview | Private preview | Private preview | Coming soon | Coming soon | Coming soon | Coming soon | Coming soon |

And we are not done. 💪

# More trigger and control-flow options

⏭ For-each loops across multiple tasks

⏭ Task groups: Visual segmentation of large DAGs

⏭ Periodic triggers: Run every week, day, hour

⏭ Multiple triggers per job

⏭ Queuing on infrastructure resources, e.g. instance pool availability

# Unity Catalog and Serverless compatibility

⏭ Compatibility checks for Unity Catalog and serverless

⏭ AI assisted code updates

# Operational health across the Lakehouse

**Single operational view**

→ Health metrics across all assets

→ Data health monitoring

→ Anomaly detection

Filter, e.g. by team or alert type

Any asset, e.g. job, pipeline, table

Any issue, e.g. data quality

**Operations monitor**

| Run as ⌄ | 2024-06-03 17:00:00 → 2024-06-05 17:00:00 📅 | Run status ⌄ | Error code ⌄ |

**Top 5 error codes** (31,353 errors)

| | |
|---|---|
| RunExecutionError | 29,876 |
| MaxConcurrentRunsExceeded | 1,009 |
| ResourceNotFound | 202 |
| DriverError | 165 |
| InvalidRunConfiguration | 101 |

Legend: 🟥 Failed  🟦 Skipped  🟩 Succeeded

| Last run | Asset | Run as | Launched | Duration | Status | Error code | Run p... | |
|---|---|---|---|---|---|---|---|---|
| Jun 05, 2024, 04:41 PM | to be scheduled (2) | 👤 maksi | By scheduler | 2s | ⊗ | DataQuality | | ⋮ |
| Jun 05, 2024, 04:41 PM | maksim-bugbash | 👤 maksi | By scheduler | 3s | ⊗ | RunExecutio... | | ⋮ |
| Jun 05, 2024, 04:41 PM | New query 04/14/... | 👤 xiang | By scheduler | 3s | ⊘ | | | ⋮ |
| Jun 05, 2024, 04:41 PM | select 1: 1 > 1 | 👤 xiang | By scheduler | 4s | ⊗ | MaxConcurr... | | ⋮ |
| Jun 05, 2024, 04:41 PM | after | 👤 bin.fe | By scheduler | 3s | ⊗ | MaxConcurr... | | ⋮ |
| Jun 05, 2024, 04:41 PM | sneaky retrieveDb... | 👤 micha | By scheduler | 30s | ⊙ | | ▪ | ⋮ |

# 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**

**Public Preview, i.e. production support**

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

**Serverless Compute**

Hands-off auto optimized compute managed by Databricks

Storage

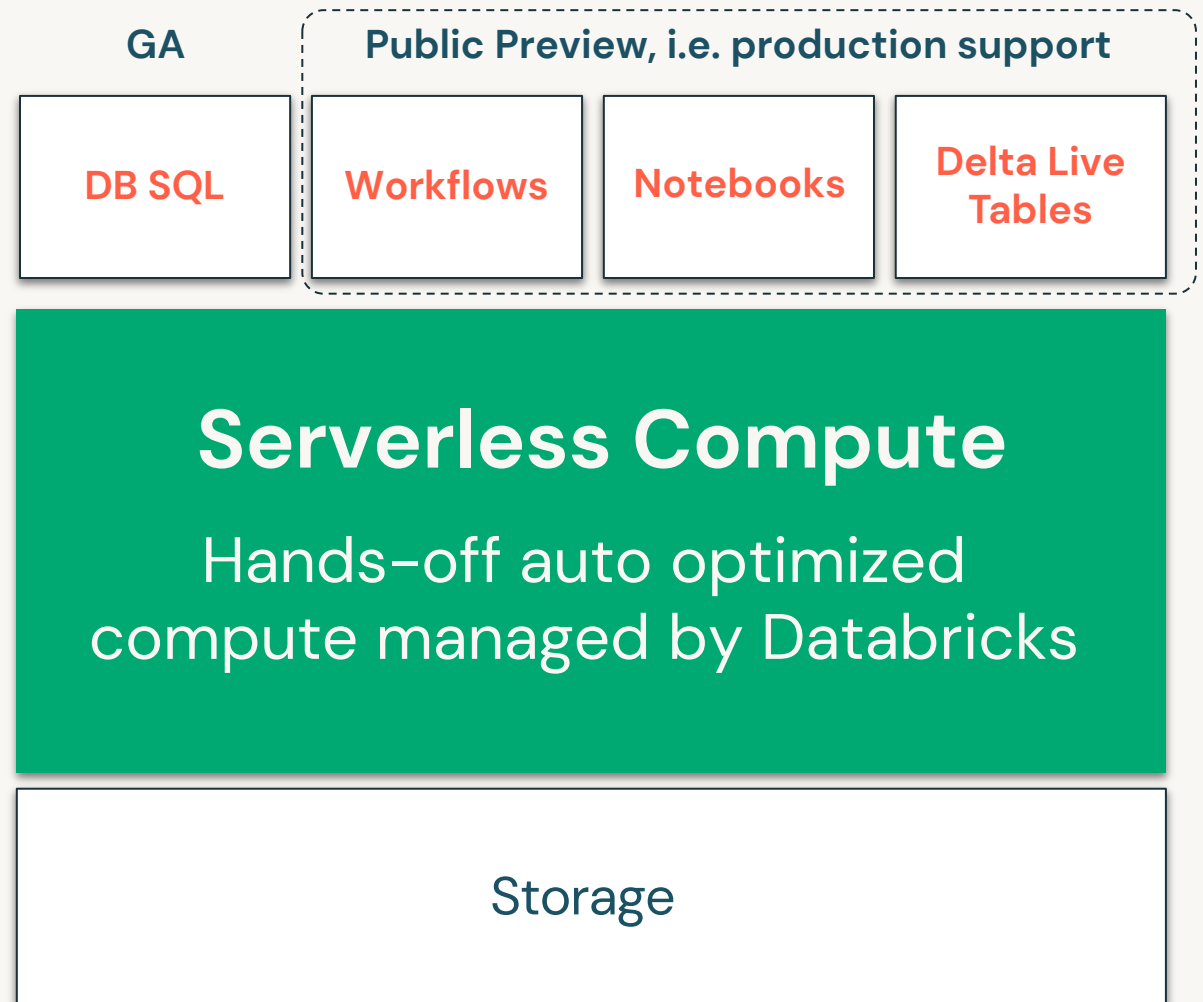Team 1    Team 2    Team 3    Team 4

VPC/VNET

VPC/VNET

VPC/VNET

VPC/VNET

**Structured, Semi-structured and Unstructured Data**
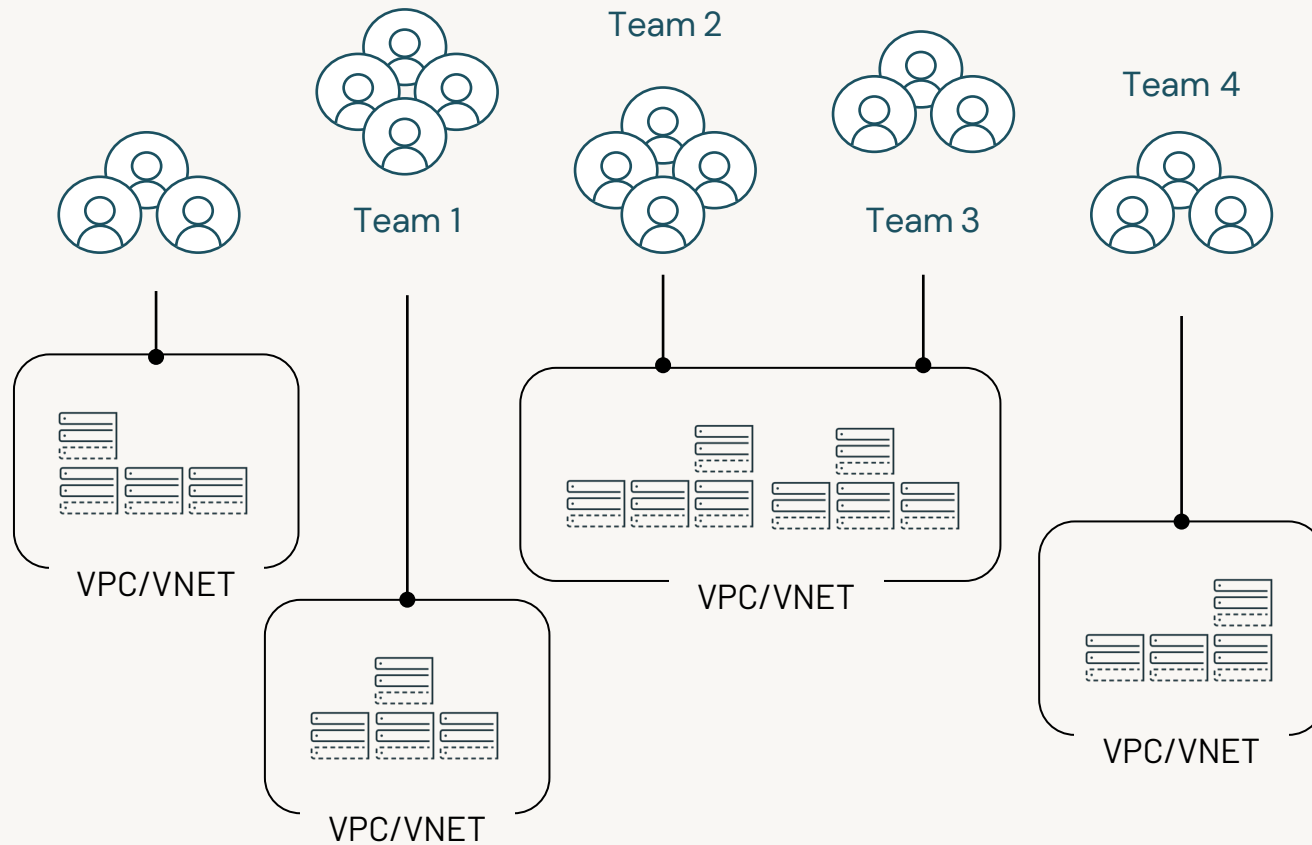
# STOP spending time on…

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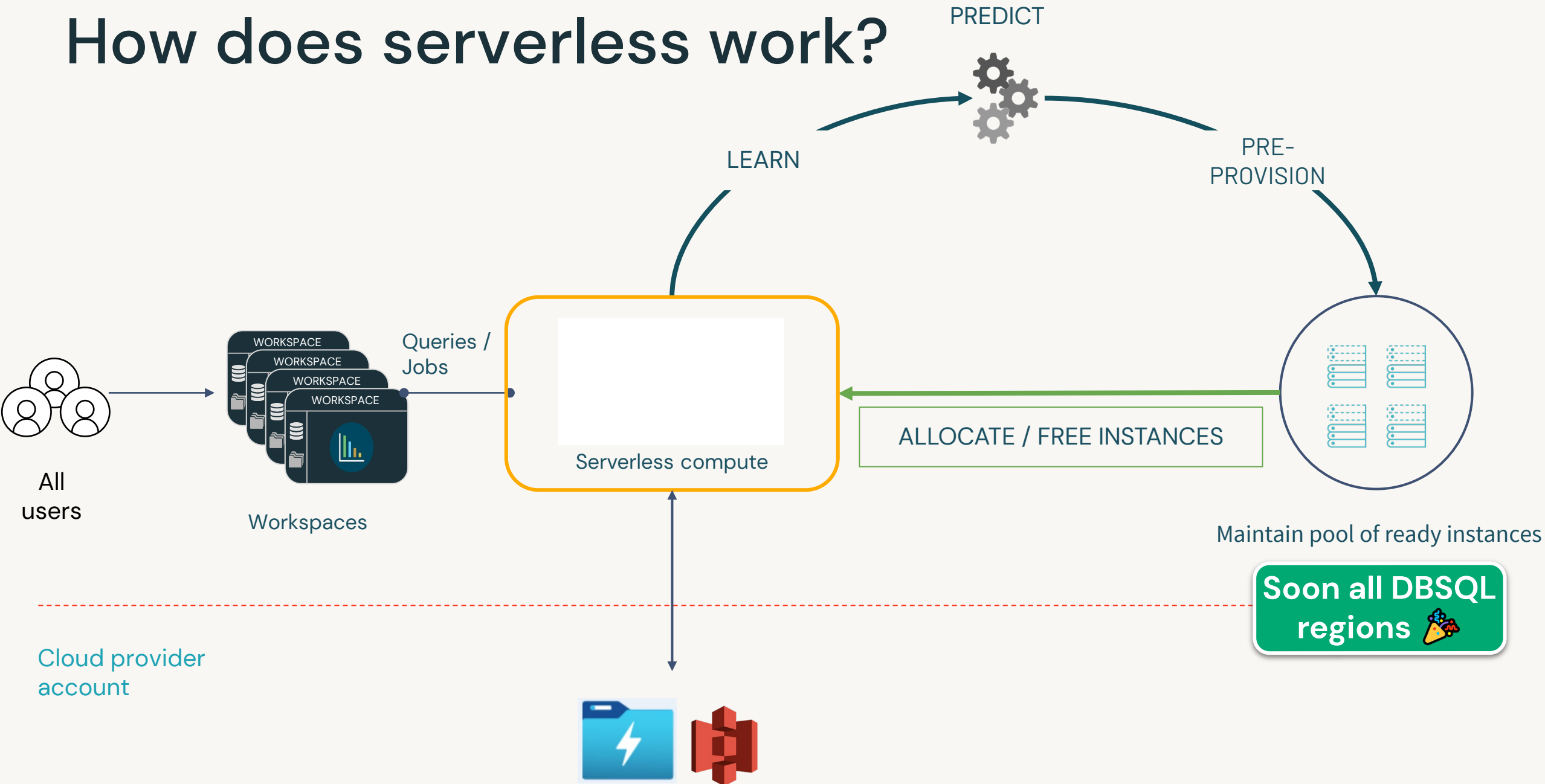
Capacity projections and reservations
Right sizing instances for workloads
Maintaining high utilization
Managing instance pools
Vacuum / compaction of Delta tables

# How does serverless work?

PREDICT

LEARN

PRE-PROVISION

Queries / Jobs

Serverless compute

ALLOCATE / FREE INSTANCES

All users

Workspaces

WORKSPACE
WORKSPACE
WORKSPACE
WORKSPACE

Maintain pool of ready instances

**Soon all DBSQL regions** 🎉

Cloud provider account

# Key technologies in Serverless

- AI managed warm-pool of VMs enabling faster up- and down-scaling

- Enhanced horizontal autoscaler

- Versionless: automatically latest features (DBR, photon, etc.)

- First and only secure multi-user Spark w/ fully isolated user code

- Environment caching

- Automatic vertical scaling (soon)

**Fully managed**

**Automatically improving**

# Serverless Compute in Workflows



**Fully managed and reliable**

- <60s startup
- Automatic failover
- Cost optimized + development mode (soon)

# When to use Serverless

| | Use cases |
|---|---|
| ▶️ **Use now** 🥳 | <ul><li>Interactive Pythons or SQL (no Scala yet)</li><li>New jobs</li><li>Existing jobs compatible with Unity Catalog shared access mode</li><li>Performance/startup time, streaming is important</li><li>Instead of instance pools or all-purpose compute</li></ul> |
| ⏭️ Later this year | <ul><li>Cost-optimized mode</li><li>Team level cost attribution</li><li>Internet access controls</li><li>GPUs</li></ul> |

- Soon all serverless SQL regions!
- Built on UC Lakeguard

# General availability of serverless compute for Notebooks, Workflows, DLT

## 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

Breakout session today at 4:00 PM

...rolling out next few weeks

...in all regions with serverless SQL

# Demo

# More to explore – Databricks Workflows

**At DAIS**

- **Practitioners Guide to Serverless Compute**
  Today, 4:00 PM PDT, South, Level 3, Rm 305
- **Ingestion Connectors**
  Thursday, 11:20 AM PDT, West, Level 2, Rm 2009
- **Workflows: Practical How-Tos and Demos**
  Thursday, 12:30 PM PDT, South, Level 3, Rm 307
- **Keynote on Data Engineering**
  Thursday, Jun 13, 8:30 AM PDT, South, Expo Level, Hall C

**...and beyond**

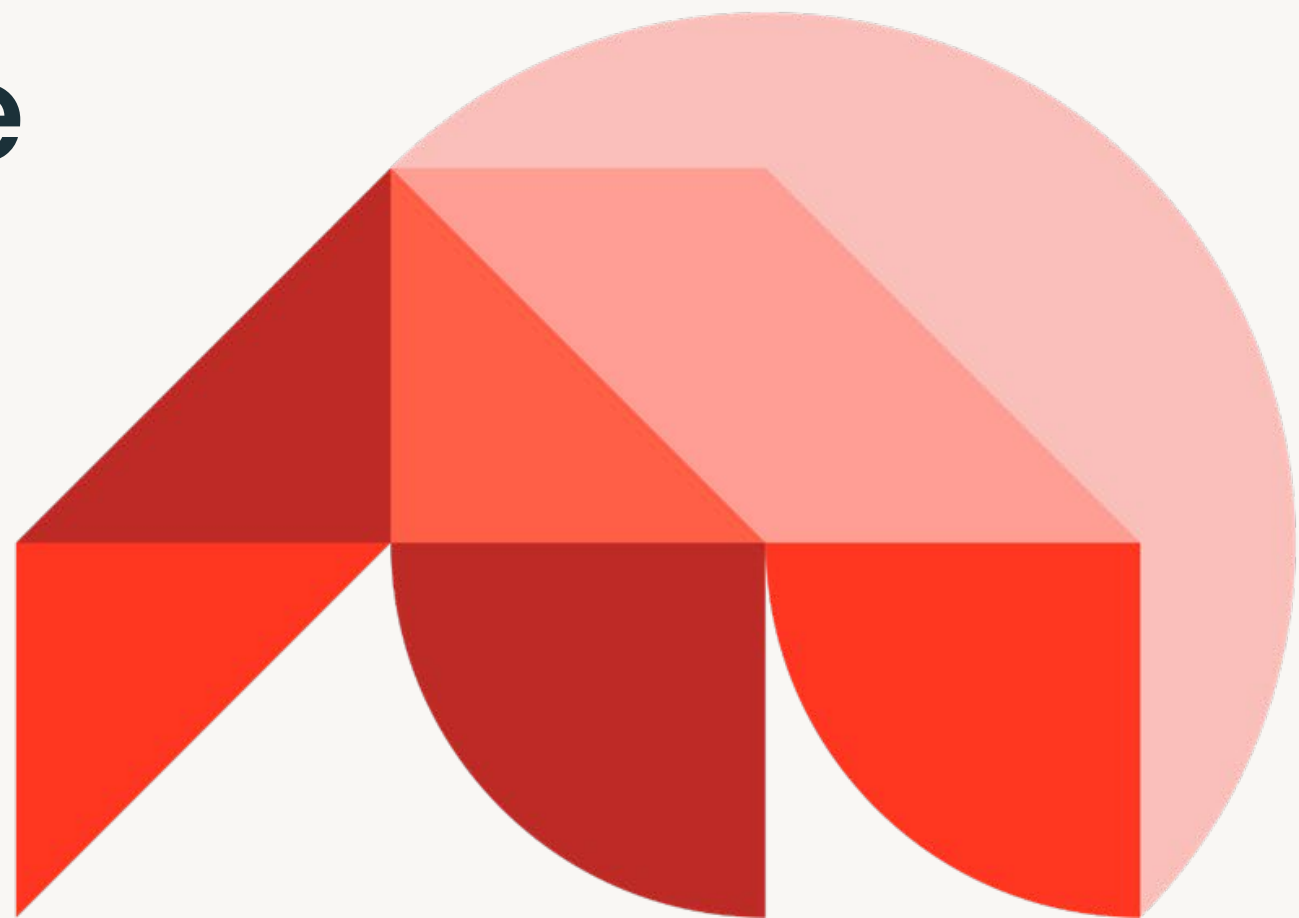- databricks.com/product/workflows
- databricks.com/demos

🙏 Thank you

# OUT

# databricks

# Your short but interesting slide title goes here
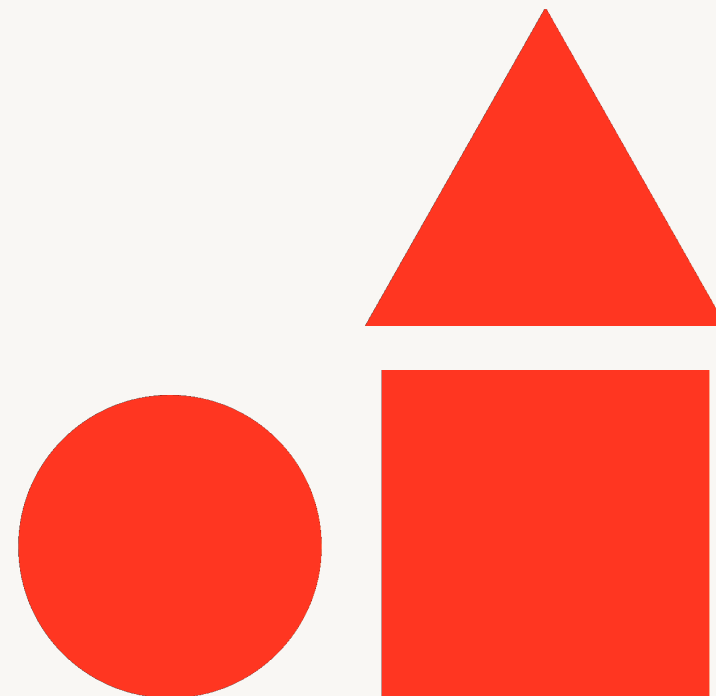
**Author Name**
Date

# Your short but interesting slide title goes here
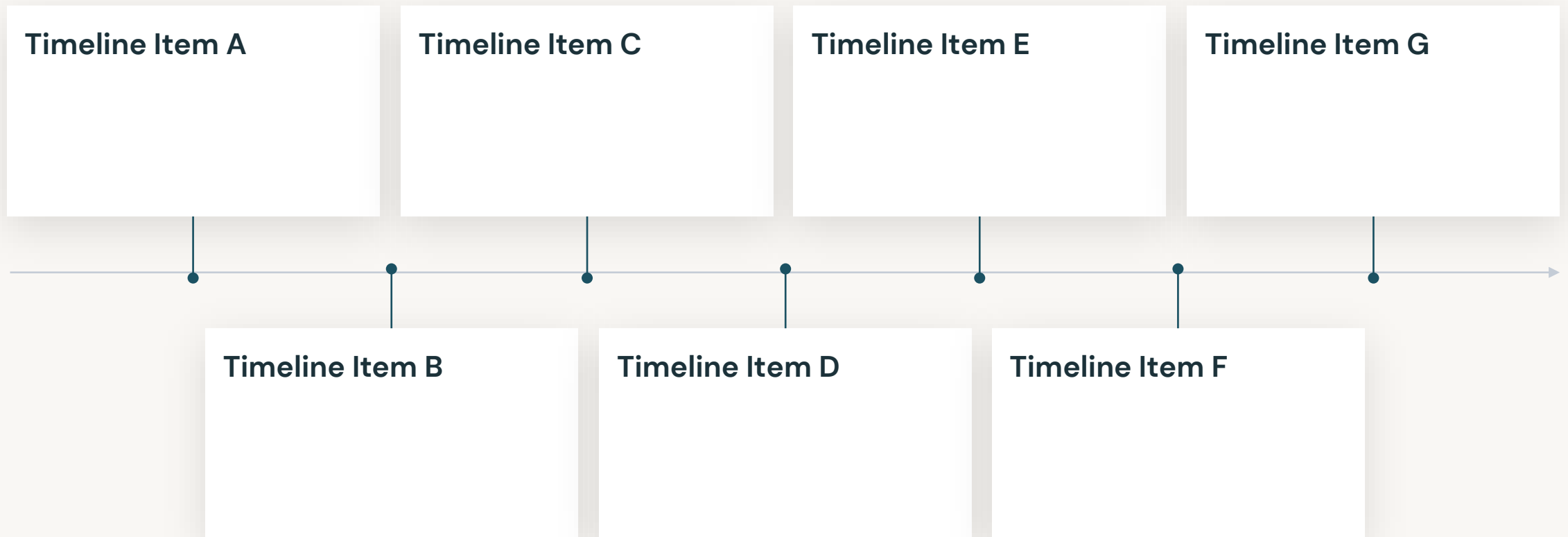
**Author Name**
Date

# Resources

# Table samples

## Your subtitle here

| Table Header 1 | Table Header 2 | Table Header 3 |
|---|---|---|
| Table Content | Table Content | Table Content |
| Table Content | Table Content | Table Content |
| Table Content | Table Content | Table Content |

| Table Header 1 | Table Header 2 | Table Header 3 | Table Header 4 | Table Header 5 | Table Header 6 |
|---|---|---|---|---|---|
| Table Content | Table Content | Table Content | Table Content | Table Content | Table Content |
| Table Content | Table Content | Table Content | Table Content | Table Content | Table Content |
| Table Content | Table Content | Table Content | Table Content | Table Content | Table Content |

# Timeline style

## Your subtitle here

**Timeline Item A**

**Timeline Item C**

**Timeline Item E**

**Timeline Item G**

**Timeline Item B**

**Timeline Item D**

**Timeline Item F**

# Primary icons

## Examples

Included are a few various icons and illustrations. To access the full library of icons, please follow this link:

**Click for primary icons**

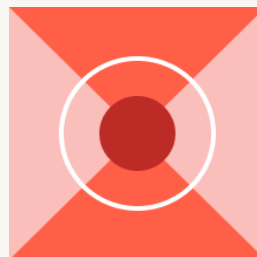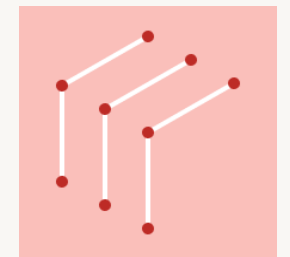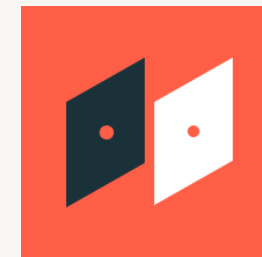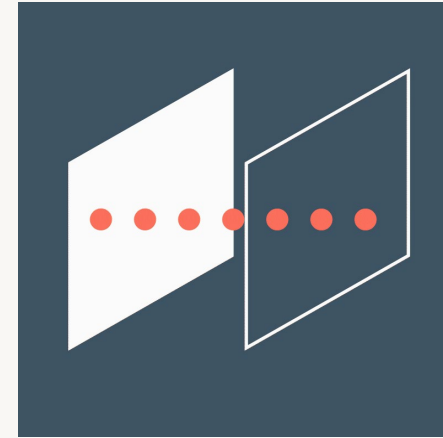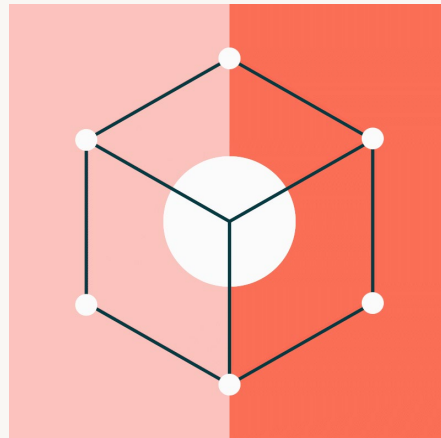| | | | |
|---|---|---|---|
| Life Sciences | Cloud Security | Analytics | Data Sharing |
| Collaboration | Retail | Multi-cloud | Gaming |
| Public Sector | Prediction | Data Science | Data Lake |

# Secondary icons

## Examples

Included are a few various icons and illustrations. To access the full library of icons, please follow this link:

**Click for secondary icons**

# Illustrations

## Examples

Included are a few various icons and illustrations. To access the full library of icons, please follow this link:
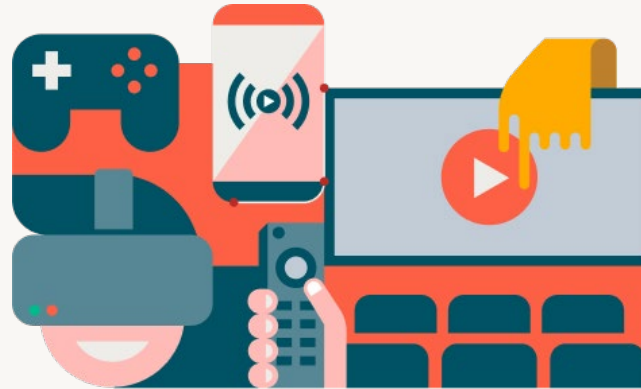
**Click for illustrations**

Manufacturing

Retail

Media & Entertainment

Healthcare and Life Sciences

Public Sector

Financial Services