# Practical Data Governance

Managing data compliance and privacy within a large scale Databricks environment

**Brad Nicholas**
Director, Digital Platforms, Corning

**Aaron Colcord**
Sr. Director, Privacera

ORGANIZED BY databricks

1

# About us

## Aaron Colcord

**Privacera**

Aaron is an Adaptive technical leader with 20+ years' experience in spearheading enterprise data solutions and enabling scalable, secure processes which lead to powerful insights from complex data systems.

He has spent the last couple of years working passionately inside evolving technologies such as Lakehouse, Data Mesh, the ever-evolving modern data stack, and acquiring 7 Patents in this area.

He joined Privacera because a belief in the mission and technology to advance Customers and their data management programs.

## Brad Nicholas

**Corning**

Brad runs the software engineering team at Corning Emerging Technology responsible for digital transformation platforms including Databricks & Privacera.

The team works exclusively with open source and open core software and is responsible for the broad-scale adoption of these technologies by digital practitioners across Corning's businesses and functions worldwide.

His experience includes software engineering, product and general management roles at multiple startups and large-scale enterprises.  He holds 12 networking and IIoT patents.

# What is practical data governance?

Realizing global data-driven scale and value

- Material financial impact through scalable value delivery, with velocity
  - This scopes everything we do
  - Global data governance, software patterns &  inner-sourcing are essential
  - Fine-grained runtime access control in the Data Lakehouse is essential as well
  - Security must be designed in, not an overlay

- Open, adaptable, enterprise-grade tech stack
  - Apache Spark, Delta Lake, MLflow
  - Apache Ranger (Data Security)

# Allow domain experts to self-serve

## Patterns, CI/CD & automated data governance are essential
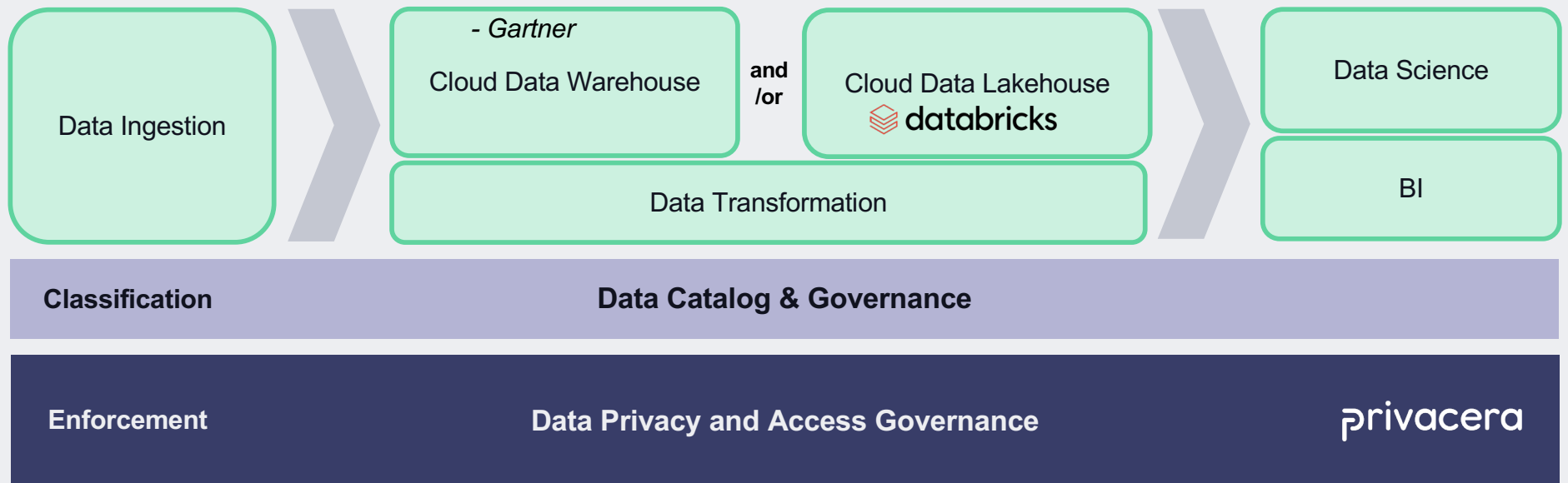
- Professional SWEs and SREs scale data-driven transformation solutions
  - Not a core competency for many successful enterprises – their differentiating expertise is elsewhere



- The key challenge is to enable domain experts to leverage advanced analytics effectively
  - Experts must focus on the right questions and have a straightforward path to data for answers
  - They act on what they learn & move on to the next question

- We're done when practitioners don't need our expertise to be productive

# Curated data accessibility empowers the data platform

## Secure landing, enrichment, serving, versioning at scale

**Data governance** is the specification of decision rights and an accountability framework to ensure the appropriate behavior in the valuation, creation, consumption and control of data and analytics.

Data Ingestion

*- Gartner*

Cloud Data Warehouse

**and /or**

Cloud Data Lakehouse
databricks

Data Transformation

Data Science

BI

**Classification** **Data Catalog & Governance**

**Enforcement** **Data Privacy and Access Governance** privacera

# Everyone's reality is existing "brownfield"

## Transformation forces a focus on the data lake / lakehouse

- Transformation data sources get processed first

- Change Data Capture decouples the source from the landing platform

- Address data governance & pipeline automation in parallel, opportunistically

- Lather, rinse, repeat

**Start early, define patterns, scale through repetition**

Time

Classify/inventory data at **early point of collection** to create downstream efficiencies

**Company data**

...but most companies classify/inventory data where **data size causes inaccuracies and delays**

# Balance governance & borderless data access

Don't compromise on either

## Drive effective data governance with Policy as Code

- Runtime data access enforcement of operational data contracts
- Agile, automated, best practice policy management
- Full compliance with legal/regulatory obligations

## Proactively leverage data access control & discovery capabilities

- Active tagging of technical metadata
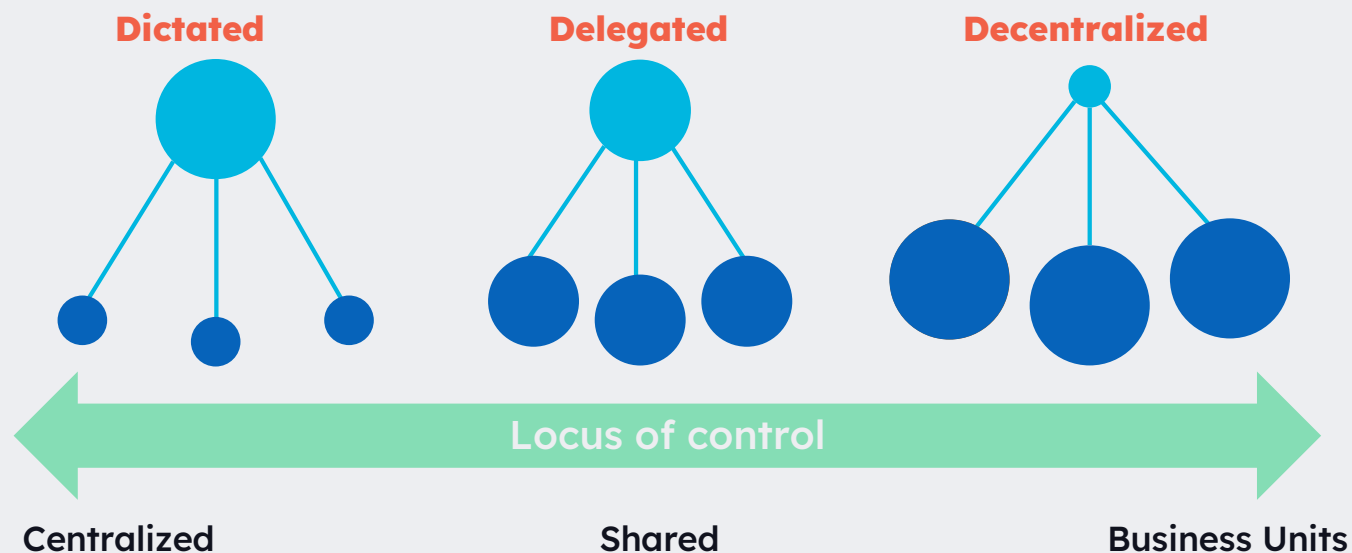- PolicyCTL as a gateway to automated version control

## Serve curated, 'Borderless' data where it enables business value

- Collapse time to data leverage
- Design apps for effective governance (e.g, avoid open fields, PII honey pots)
- Enable Virtual Boundaries that can shift and adapt

# Democratization vs Governance

Avoid sacrificing productivity

# Common pitfalls of data governance
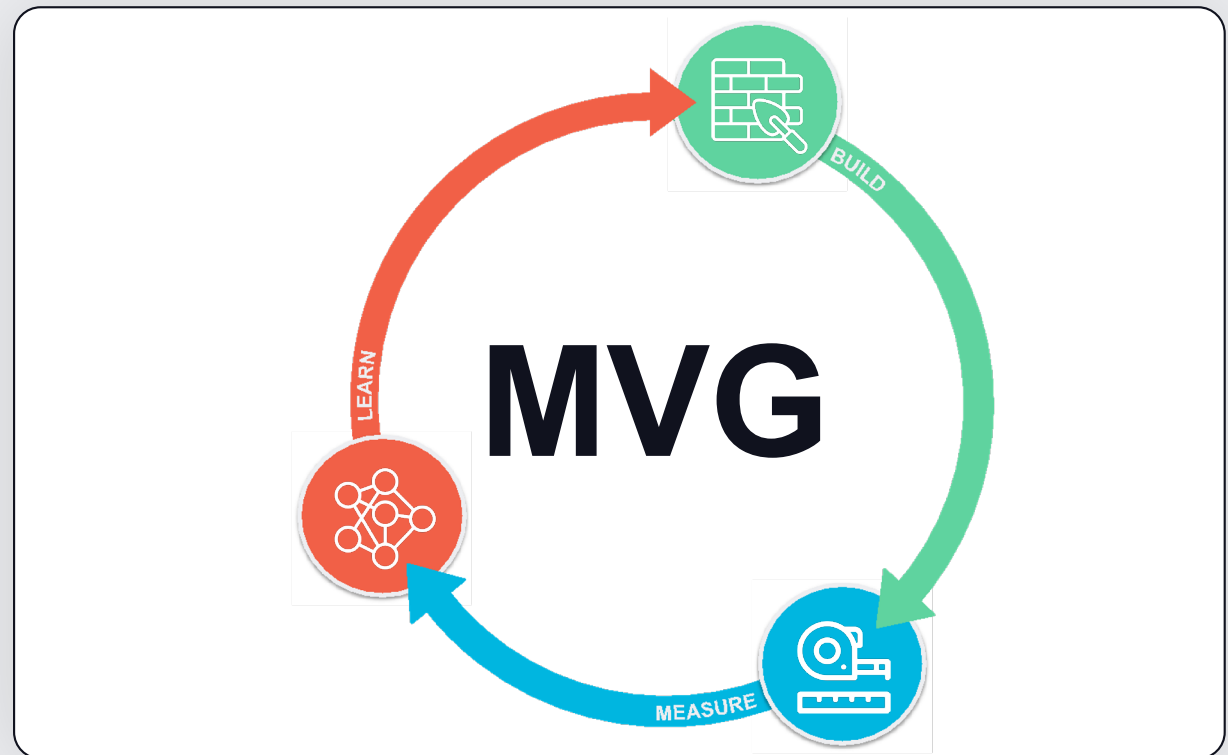
5 Challenges that lead to Failure

- **Connecting data governance to business outcomes**

- Difficulty defining and adopting consistent data governance processes and policies

- Failure to define common enterprise data definitions

- **Inconsistent approach to data across projects**

- Difficulty defining and sustaining a path to target state (competency to implement)

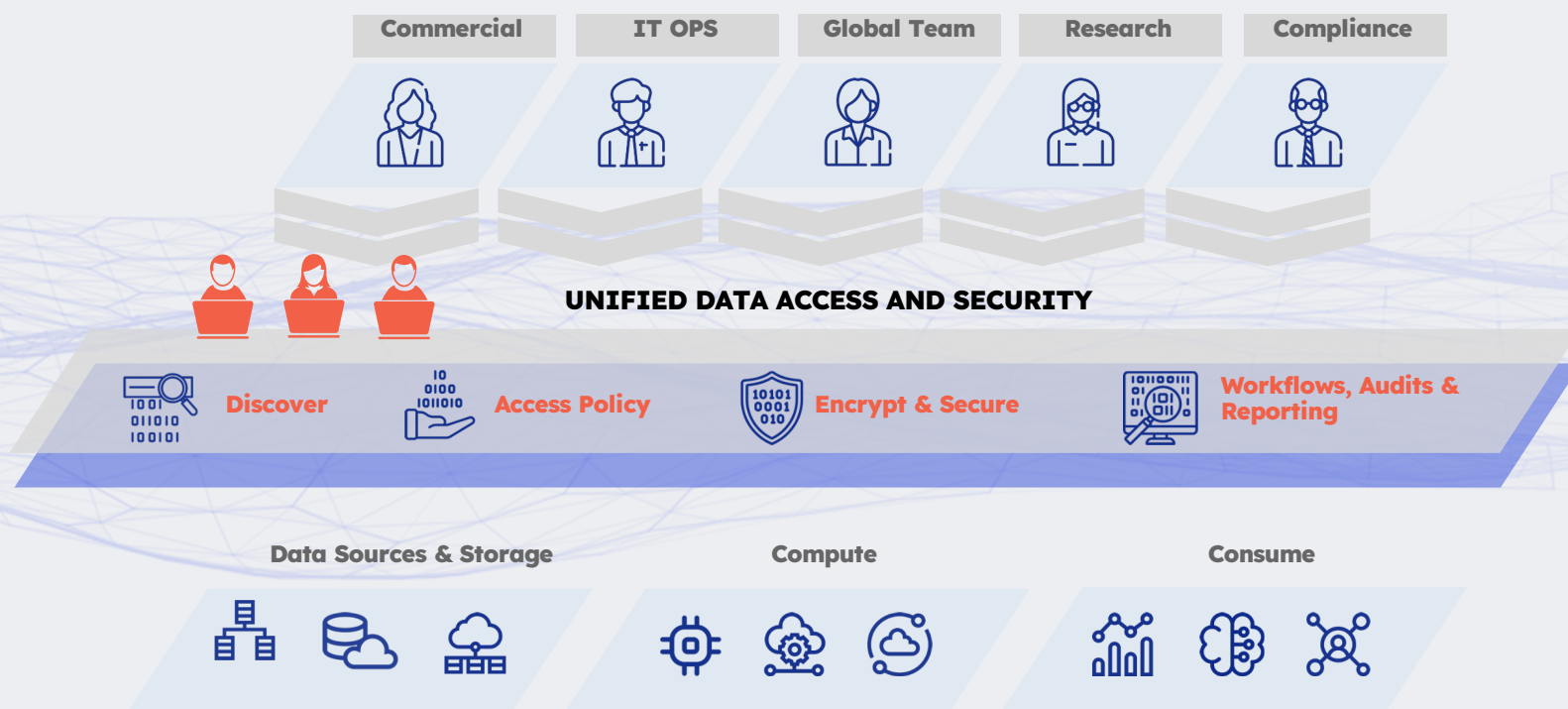# Minimal Viable Data Governance

**The paradox is real**

**Business value is…**

- Data contract is the spirit of MVG

- Most 'data products' have no defined value.

- Without quantified value, what is the justification?

- 'Rogue' business groups

- Usually the focus is on the symptoms, not the cure.

# Unified Data Access Governance

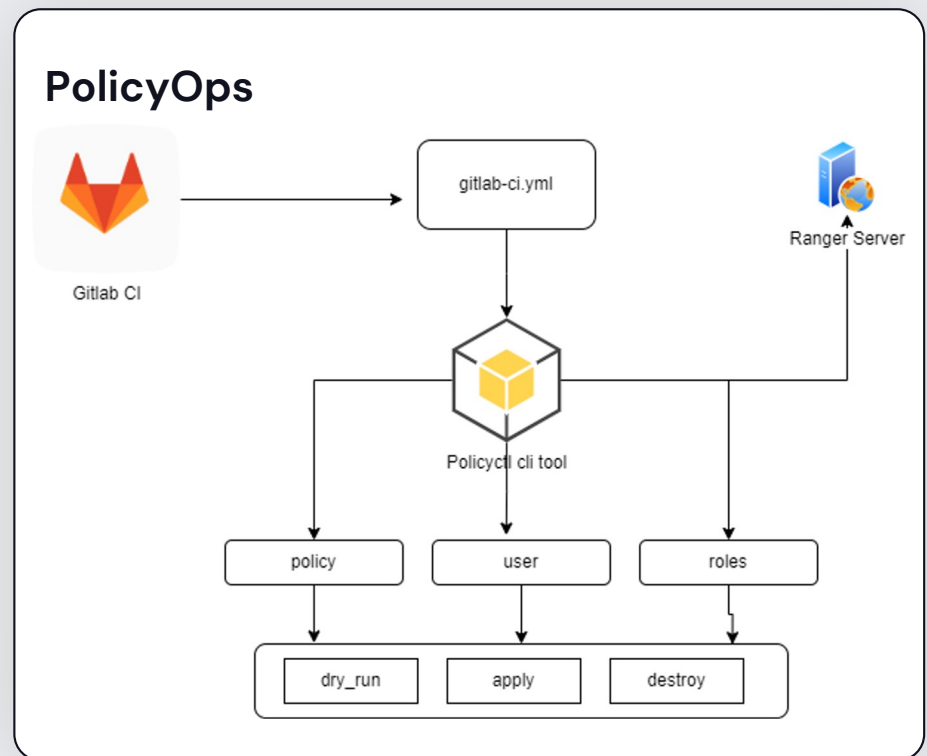## A single location for Data Access and Security Governance

| Commercial | IT OPS | Global Team | Research | Compliance |
|---|---|---|---|---|

**UNIFIED DATA ACCESS AND SECURITY**

Discover    Access Policy    Encrypt & Secure    Workflows, Audits & Reporting

Data Sources & Storage          Compute          Consume

# What is PolicyOps?

## Automating for Scale

### Open Content

Policy Ops is CI/CD tool for Apache Ranger policy management. It provide flexibility to store policy, user, roles etc., inside version control so any changes will be controlled using git and it will be easy to maintain and deploy. Policy ops is based on a cli tool called policyctl which is developed by **Privacera** for its customer to make changes in the Apache Ranger policies using cli commands.

- YAML based properties for easy understanding
- More Controlled and Error Free Changes.
- Easy command-line step execution.

## PolicyOps



Gitlab CI → gitlab-ci.yml

Ranger Server

Policyctl cli tool

| policy | user | roles |

| dry_run | apply | destroy |

# Scaling techniques

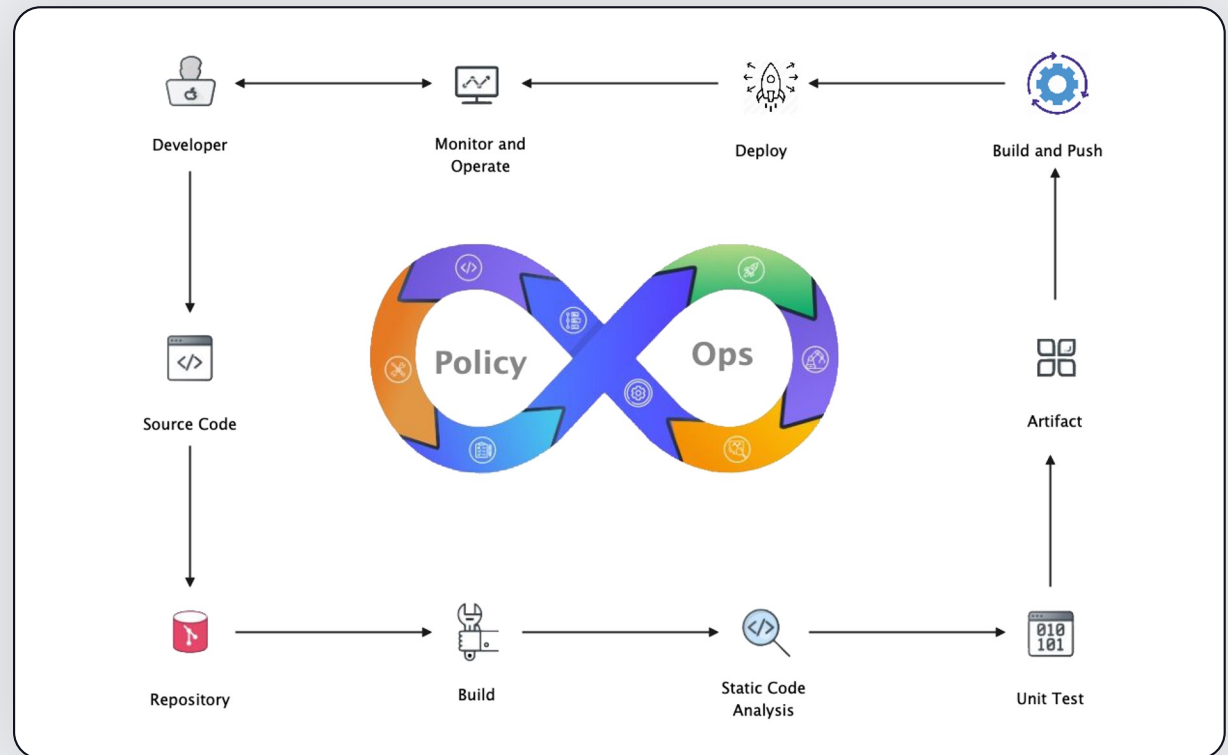Automating access control through PolicyOps

**Policy as Code**

- Roles
- Data contracts
- Test datasets

**Resources & Access Methods**

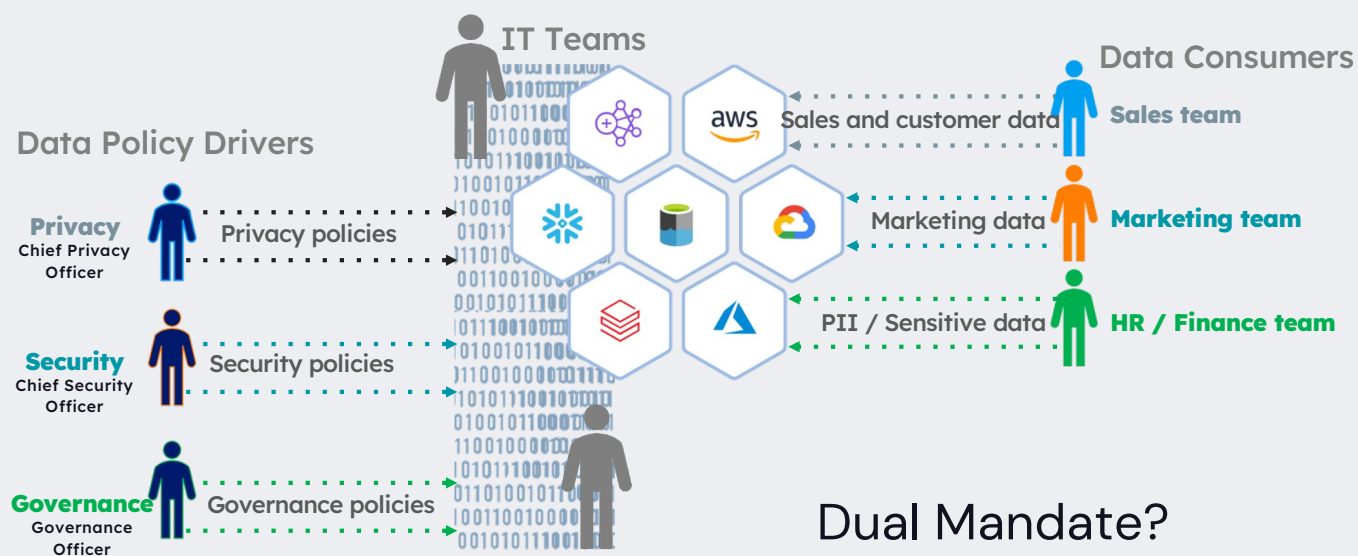- Notebooks
- Jobs
- Files
- Buckets

**Terraform**

- Plan/Apply is your friend

DATA+AI
SUMMIT 2022

# Our demo scenario

One of our enterprise stakeholders now has a concern



IT Teams

Data Consumers

Data Policy Drivers

Sales and customer data — Sales team

Privacy
Chief Privacy Officer — Privacy policies

Marketing data — Marketing team

Security
Chief Security Officer — Security policies

PII / Sensitive data — HR / Finance team

Governance
Governance Officer — Governance policies

Dual Mandate?

**Former AWS engineer convicted over hack that cost Capital One $270m**

# Solution Ingredients

**Encryption Scheme:** Defines how to encrypt the Data. What users can decrypt the data

**Presentation Scheme:** Can Keep data sensitive hidden based on User
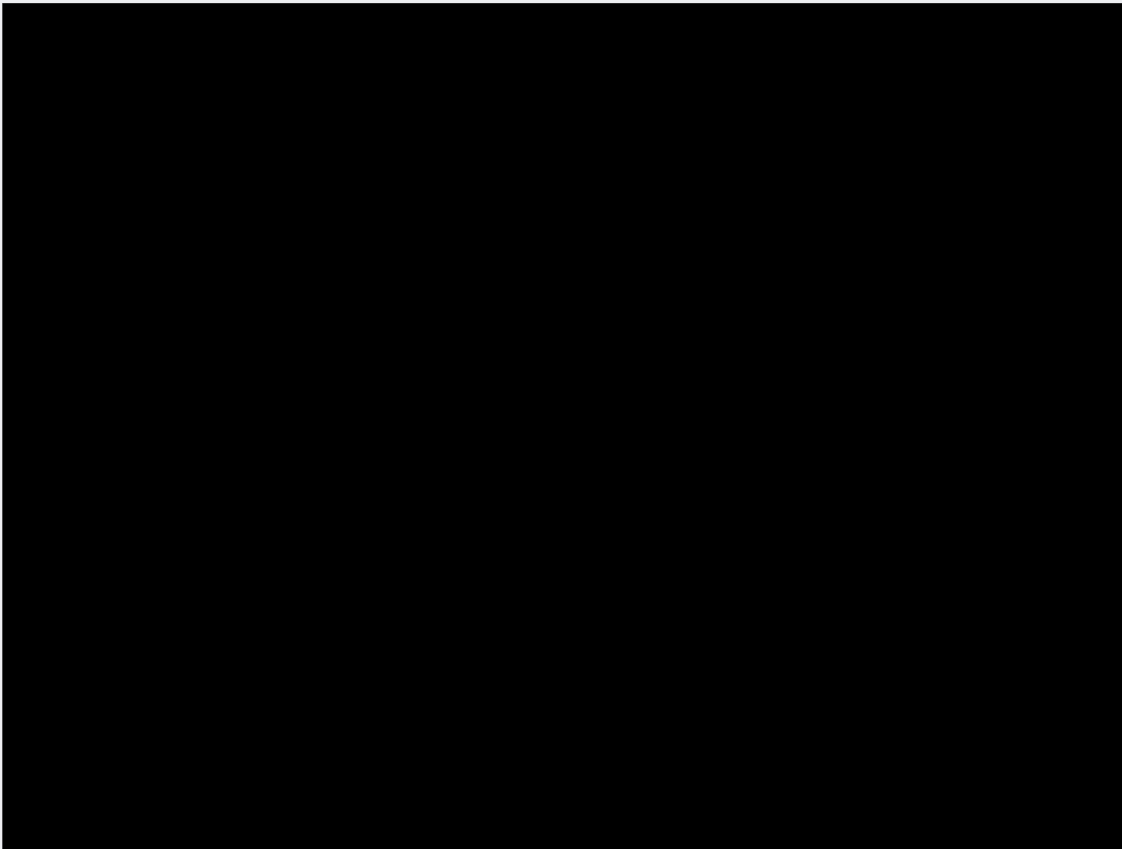
Keys can be stored externally

Encryption can occur automatically in a workflow
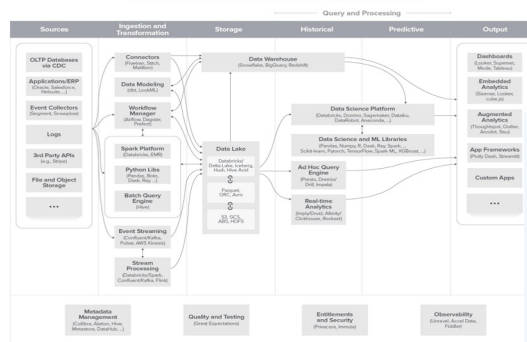
# Demo

PolicyOps

# Demo

## PolicyOps

```
1
2
3    Clone existing ranger server policy , user , role etc. this will download existing policies from the server and put it inside a git folder.

     policyctl package create --name path/to/my/package
4
     Test what are the changes made inside yaml file . We run command in –dry_run mode which  don't apply the changes but give a glimpse of
5    the changes.

6    policyctl package apply --name path/to/my/package –dry_run=TRUE

7    Apply changes for a single policy.

8    policyctl policy apply --name path/to/my/package/policy.yaml

9
10
```

# Open-source
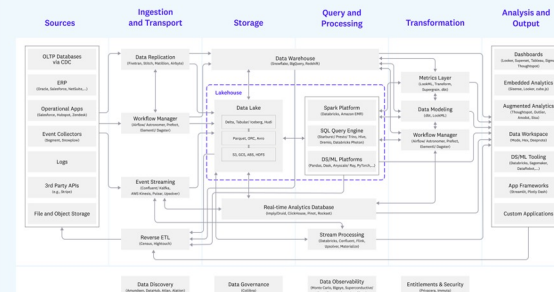
- Ability to adapt to your business

- Scalable

- Interoperable

- No Vendor Lock-In



**Rapidly adapt and Evolve**

Apache Ranger