

# **Scaling Privacy**

## Practical Architectures and Experiences



Aaron Colcord Sr. Director, Center of Excellence, Privacera

1



**Mei Gui** Software Engineer, Databricks

ORGANIZED BY Sdatabricks

# 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 of a belief in the mission and technology to advance customers and their data management programs.

## Mei Gui

#### Databricks

Mei is a software engineer working on DBFS Fuse and files in repos at Databricks. Mei is the lead engineer for integrating Privacera and DBFS Fuse on Databricks. Prior to Databricks, Mei was a software engineer at Uber. "If we accept as normal and unavoidable that everything in our lives can be aggregated and sold, then we lose so much more than data. We lose the freedom to be human."



Tim Cook CEO, Apple



# Review from last year's presentation

## The Main Points

## Have you ever...

Gone to a website and read their privacy policy, clicked accept cookies, accepted terms of service, or EULA?

#### Collecting information about your customers can

- Improve the experience.
- Allow the company to understand their business better.

#### At the core, privacy is a policy and legal obligation

- You have the data. It used to be your business to just secure it.
- Do you want your information monetized? Sold? Traded?
- Most companies don't do this. But the privacy policy is there for you.

#### Clicked 'accept all' on website, used a digital assistant..

# Review from last year's presentation

#### The Main Points

## What is the difference between these?

#### Security

- Preventing unauthorized usage of systems
- Ensuring users don't see the incorrect information
- Creating boundaries to enforce the right action in the system

#### Compliance

The process of making sure your company and employees follow all laws, regulations, standards, and ethical practices that apply to your organization

#### Privacy

- Data privacy may be defined as the authorized, fair, and legitimate processing of personal information
- Consent rights
- Do not share
- Slippery space

# Review from last year's presentation

### The Main Points

## Ideal Scalable System



# **Practical Architectures**

# **Project Mercury**



## Data eco-systems

## **Evolving Architectures**



# Our story

## **Being Democratized**



# Privacy Impact Assessment (PIA)

## Project by Project Process

Define the Scope

#### Data Flow Analysis -

- Show the lineage
- Identify the PII
  - □ Why is it being collected (Business Reason)

Privacy Analysis - How are we using this data and for what purpose?

**Report - Risks and Implications** 

## Demo

## Governed Agile Data Science

	N	Lflow Training a	nd Inference	with Priva	Acera Python						∅ s	chedule ~ Share
0-	4	privacera-demo	I∨ 🖹 File ▼	💣 Edit 🕶	🌆 View: Standard 🕶	🕑 Run All	🖉 Clear 🕶	O Help		🙊 Comments	A Experiment	D Revision history
	a	Cmd 1										
•		MLflow q	uicksta	rt: trai	ining and	loggi	ng					
જ	This tutorial is based on the MLflow ElasticNet Diabetes example. It illustrates how to use MLflow to track the model training process, including logging model parameters, met model itself, and other artifacts like plots. It also includes instructions for viewing the logged results in the MLflow tracking UI.											s, metrics, the
0												
Q		This notebook uses the scikit-learn diabetes dataset and predicts the progression metric (a quantitative measure of disease progression after one year) based on BMI, blood pressure, and other measurements. It uses the scikit-learn ElasticNet linear regression model, varying the alpha and ll_ratio parameters for tuning. For more information on ElasticNet, refer										
~		<ul> <li>Elastic net regular</li> </ul>	ization									
æ		Regularization and	d Variable Selection	n via the Elas	tic Net							
ஊ			_									
		Requirements										
	<ul> <li>This notebook requires Databricks Runtime 6.4 or above, or Databricks Runtime 6.4 ML or above. You can also use a Python 3 cluster running Databricks Runtime 5.5 LTS ML.</li> <li>If you are using a cluster running Databricks Runtime, you must install MLflow. See "Install a library on a cluster" (AWS Azure GCP). Select Library Source PyPI and e the Package field.</li> </ul>										cks Runtime 5.5	LTS or
窃											nter mlflow in	
0		<ul> <li>If you are using a</li> </ul>	cluster running Dat	abricks Runt	time ML, MLflow is al	ready installe	d.					
۲		Note										
٨		This notebook expects that you use a Databricks hosted MLflow tracking server. If you would like to preview the Databricks MLflow tracking server, contact your Databricks sales representative to request access. To set up your own tracking server, see the instructions in ML flow Tracking Servers and configure your connection to your tracking server by rupping										ks sales
		mlflow.set_tracking_u	iri.	oogle Chrome	, dasking server, see	are mondette		Hueking oct ver	o una configure your		Sar tracking serv	or by running

# Data Governance

## Blink and you missed it.



# Did we pass the PIA?

#### Did we balance Democratized with Governance? What we saw in action

- Fine-Grained Access Control- All Data, Models, Logs were controlled by One Policy (Delta, S3, and MIFlow were managed in spot)
- Audit trails showing exactly what was working
- Unstructured Data governed in Development

#### What we had happening in the background

- Data Movement was being tracked
- What type of Data (PHI) was being used.

#### Building this represents true Business Value

- Data Lakes power this via 'Borderless Data'
- We 'Scaled Privacy'
- The PIA was a living, breathing enabling element; not a static process requiring high friction coordination.



DA'IA+AI SUMMIT 2022

# **Practical Experiences**

# How do you take this and build a program?

# Set some principles for architecture

## Make your decisions from the principles

- Build around Services, not libraries.
  - Avoid 1 Million Hour Problems
- Implement Policy-Approaches
- Separation of Duties
- Be Actionable
- Plan to be Scalable from the start.
- Thread the *Build vs Buy* Equation

# Authorization

Authorization is hard. Authentication is not.

Lean into policy-based controls that allow:

- Source access control with the raw data
- Managed control that enforces how data is used
- Rapid updates while the true definition of 'sensitive information' is found
- Enablement from the ability to adapt to regulations
- Quality Control to influence



# Organizing your Lake, Warehouse, Cloud

### **Practical Experience**

#### We want to think in terms of the business

- Example: Customer, Field, Product or Customer, Flight, Marketing
- If we align this way, at least our ontology will build the shape

#### Organize Objects close to these terms and have some way to map

- Losing track of S3 buckets is a problem. Yet it's their flexibility that makes them useful.
- You can't find it later if you didn't plan to find it later.

#### **Tagging Strategy**

- Pieces of data may be related to larger subjects
  - Example: Who took the flight is a Customer and stored as a part of Flight.
- We can attach information flexibility about origin, PII, Technical Context, Business Context

# How have you managed this?

## Engineers vs. the world

#### All of these processes exist

GSA - A lot of great material

- PIA
- Procedures
- Training Requirements

#### **HIPAA Directives**

- How to de-identify
- Safe Harbor
- Expert Determination

#### "Build it and they will come"

Unify your Stakeholders

- Most of these tools are DIY
- The moment you make the tools, you need adaptability
- Open Source bridges the 'Buy/Build'
- Policy-Approach can solve by connecting Tags to 'Intent to Use'
  - Use Masking Policies to change how data is returned dynamically
  - ABAC

# Separation of duties

Don't let the mouse guard the cheese!

#### This isn't that hard.

- It doesn't always mean people; it can be process.
- Your expert to help defend the organization can't be the requestor



## DATA+AI SUMMIT 2022

# Thank you