

### Distributing Data Governance: How Unity Catalog Allows For A Collaborative Approach

Gilad Asulin Pulkit Chadha

June 27th 2023



### **Your Presenters**



#### **Gilad Asulin**



**Akamai** Sr. Big Data Team Lead @ Akamai





#### **Pulkit Chadha**



Sr. Solutions Architect @ Databricks



@pulkitchadha

### What You will Learn Today?



Why Unity Catalog is essential for Databricks users?



How can you start your journey to use Unity Catalog?



What utilities can be built using Unity Catalog's metadata?

## Akamai Technologies Inc

The world's most distributed platform for cloud computing, security, and content delivery



**Everywhere** you do business, and **anywhere** customers come online, Akamai is closer, with more capacity and integrated security, acceleration, and computing capability than anyone.

# Akamai Application Security (AppSec)



### Akamai AppSec Lakehouse



@giladasulin, @pulkitchadha

#### **In Numbers**



#### **In Numbers**



#### **In Numbers**



#### **In Numbers**



### Why Do We Need Unity Catalog?

#### What Unity Catalog Helps With?



What Unity Catalog Helps With?

Data Sharing

**Behavior:** Can't share managed tables across workspaces

What Unity Catalog Helps With?

Data Sharing

**Behavior:** Can't share managed tables across workspaces

**Optional Approach**: Create external tables

- A lot of boilerplate code
- Mount points are unsecure & unreliable

What Unity Catalog Helps With?

#### Data Sharing

Least Privileges

**Behavior:** Can't share managed tables across workspaces **Behavior:** There is no option to restrict users in workspace level

**Optional Approach**: Create external tables

- A lot of boilerplate code
- Mount points are unsecure & unreliable

### What Unity Catalog Helps With?

#### Data Sharing

**Behavior:** Can't share managed tables across workspaces

**Optional Approach**: Create external tables

- A lot of boilerplate code
- Mount points are unsecure & unreliable

Least Privileges

**Behavior:** There is no option to restrict users in workspace level

**Optional Approach**: Create separate workspaces

- Leads to difficulty in data sharing
- Causes duplicate work between workspaces

### What Unity Catalog Helps With?

#### Data Sharing

**Behavior:** Can't share managed tables across workspaces

**Optional Approach**: Create external tables

- A lot of boilerplate code
- Mount points are unsecure & unreliable

Least Privileges

**Behavior:** There is no option to restrict users in workspace level

#### **User Governance**

Behavior: Standalone user management system (w/o using AAD)

**Optional Approach**: Create separate workspaces

- Leads to difficulty in data sharing
- Causes duplicate work between
   workspaces

@giladasulin, @pulkitchadha

### What Unity Catalog Helps With?

#### Data Sharing

**Behavior:** Can't share managed tables across workspaces

**Optional Approach**: Create external tables

- A lot of boilerplate code
- Mount points are unsecure & unreliable

#### Least Privileges

**Behavior:** There is no option to restrict users in workspace level

#### **User Governance**

Behavior: Standalone user management system (w/o using AAD)

**Optional Approach**: Create separate workspaces

- Leads to difficulty in data sharing
- Causes duplicate work between
   workspaces
- Duplicate user management for each workspace
- Unable to govern user access control

# Introduction to Unity Catalog



# Databricks Unity Catalog

Unified governance for data, analytics and AI

- Unified visibility into data and AI
- Single permission model for data and AI
- Al-driven monitoring and reporting
- Open data sharing



### **Unity Catalog -Components**



### **Unity Catalog External Locations**

Improves data governance and simplifies data access administration



# Centralized Metadata, Identity, and Access Management

#### Without Unity Catalog



#### With Unity Catalog



#### \* Important: only <u>one</u> UC metastore per region!

22 @giladasulin, @pulkitchadha

### Life of a query without Unity Catalog



# Life of a query with Unity Catalog



## **Migrating Workloads**



**<SA>** System Account (Service Principal or Managed Identity for Azure, IAM Role for AWS, Service Account for GCP)

#### 25 @giladasulin, @pulkitchadha

### All Purpose clusters with User Isolation mode

- Use User isolation mode for general workloads (ETL, data exploration, ...) using SQL and Python
- Multiple users can work on the same cluster

#### All Purpose clusters with Single User mode

- Use Single User mode for Scala users and for Data Scientists
- ML Runtime with MLflow is supported
- Only the owner can execute code on this cluster, so for notebook collaboration, co-workers can see everything, but they cannot execute cells.

#### **SQL Warehouse**

• Use SQL Warehouses for Business Analysts either using Databricks SQL Editor or external BI tools like Power BI, Tableau, ...

# Unity Catalog In Action



### Things We Took into Consideration

#### Inquiry Leads to Requirements



### Things We Took into Consideration

Inquiry Leads to Requirements



### Things We Took into Consideration

Inquiry Leads to Requirements





Governance Model Requirements	Groups > Add group > Add group
<ul> <li>Product owners responsible for their own data</li> <li>Govern w/o being a bottleneck</li> </ul>	Add group Create a group to manage users in this account. Users and role * Group name
<ul> <li>Design Decisions</li> <li>Apply a distributed data governance</li> <li>Each catalog has a product owner group</li></ul>	product1_owners       Save       Cancel

#### 31 @giladasulin, @pulkitchadha



#### 32 @giladasulin, @pulkitchadha

### Managed Location for Catalogs & Schemas

Differentiate between products and between prod/non-prod environments

Create a new external location × An external location is a cloud storage url (and paired credential) that allows access to data	A schema is the second layer of Unity Catalog's three-level namespace and organizes tables and views. Learn more Schema name test
stored on your cloud tenant. Learn more Copy from mount point ~ External location name	Storage location (optional)         Image location         Image location
bot-catalog-mng-location       URL       abfss://manage-location@botcatalogmnglocation.dfs.core.windows.net/	Create a new catalog × bles will be stored. If not specified, the location will A catalog is the first layer of Unity Catalog's three-level namespace and is used to organize your data assets. Learn more Catalog name
Storage credential Or bot-catalog-mnglocation	Storage location (optional) Cancel Create
Comment (optional)	abfss://manage-location@botcatalogmnglocation.dfs.core.windows.net         Location in cloud storage where data for managed tables will be stored. If not specified, the location will default to the metastore root location.
Read only	Comment (optional)
	Cancel Create

#### Data Engineering Team's Role

#### Requirements

- Ownership of all shared data sources
- Products' data is governed by high level standards and utilities
- Enable access to metadata

#### **Design Decisions**

- Create a common\_data catalog
- Create utility that expose metadata
- Create dashboards and alerts to gain a high level perspective

Data		
→ 🗌 commo	n_data	

### **Stakeholders Impact**

#### Requirement

 Identify and map all the potential impacts on stakeholders

#### **Design Decisions**

- Document all the scenarios that can impact the roles in the organization
- Communicate with stakeholders to make them engaged



Utilities Built using Unity Catalog's Metadata



### **SYSTEM Catalog**

#### Unity Catalog's Metadata

- Stores all Unity Catalog's metadata and Databricks workspaces activities
- INFORMATION\_SCHEMA stores metastore metadata.
  - Hive\_metastore is not part of metastore
- ACCESS SCHEMA stores the account level activities
- BILLING SCHEMA stores the billing logs in the account



### **Expand Data Assets Visibility**

#### Catalogs as PDF - Utility for non Unity Catalog users

- Contains all the catalogs\schemas\tables information
  - Comments\tags on each entity.
  - Who owns this catalog?
  - What entities we have?
- The audience
  - Employees that don't have permissions to UC like
    - External departments
    - Product managers
    - Restricted users

#### CATALOG: cost\_analysis

Cost Analysis based on WAF DLR sample 1% Owner: data\_engineering

#### SCHEMAS:

#### Schema cost\_analysis.prod

\*\*Production Schema for Cost Analysis\*\* Owner: data\_engineering

#### cost\_analysis.prod.actions

A statistics table showing the client's actions in response to suspicious events  $\mathsf{Owner:}\xspace$  data\_engineering

Column Name	Comment				
applied_action	The applied action by the customer				
count	how many time the action applied by customers				
count_pcnt	The percentage for this action in the specified date				
report day	the report day in dd/MM/YYY				

### **High Level Governance Dashboards**

#### State of the Lakehouse



#### @giladasulin, @pulkitchadha

### High Level Governance Dashboards

#### **Users' Activity Monitor**

#### Monitoring User Activities + Add tag

Deletion Activity								Last Week Permissions Gr	anted				
email	version created_a	t	created_on 🤤	workspace_i	d	source_ip_address	s u	grantor	grantee	table_catalog	table_schema	table_name	
in@akamai.com	2.0 2023-0	6-05 18:10:43.048	2023-06-05	7799	33222362	19	N	/@akamai.com	v@akamai.com	csi_common_data	information_schema	views	
*muner@akamai.com	2.0 2023-0	6-01 16:51:51.001	2023-06-01	7799	33222362		N	/@akamai.com	data_engineering	csi_common_data	information_schema	views	
er@akamai.com	2.0 2023-0	6-01 13:04:32.025	2023-06-01	7799	33222362		N	er@akamai.com	data_engineering	dataeng	hackathon	query_costs	
er@akamai.com	2.0 2023-0	6-01 12:03:09.702	2023-06-01	7799	33222362		N	er@akamai.com	data_engineering	apr	information_schema	table_privileges	
er@akamai.com	2.0 2023-0	5-30 19:05:52.013	2023-05-30	7799	33222362		N	er@akamai.com	data_engineering	apr	information_schema	schemata	
er@akamai.com	2.0 2023-0	5-30 17:52:45.859	2023-05-30	7799	33222362		N	er@akamai.com	v@akamai.com	cost_analysis	test	view_weekly_top_acounts	_stats
n@akamai.com	2.0 2023-0	5-30 11:43:52.729	2023-05-30	5981	15293997		A	er@akamai.com	v@akamai.com	cost_analysis	test	dlr_count_from_file	
n@akamai.com	2.0 2023-0	5-30 11:42:51.693	2023-05-30	5981	15293997		A	er@akamai.com	v@akamai.com	cost_analysis	test	view_top_10_bvm_action:	s
n@akamai.com	2.0 2023-0	5-30 11:42:21.478	2023-05-30	5981	15293997		A	er@akamai.com	data_engineering	apr	information_schema	schema_privileges	
n@akamai.com	2.0 2023-0	5-17 14:07:42.980	2023-05-17	5981	15293997		A	er@akamai.com	/@akamai.com	cost_analysis	test	view_bvm_weekly_average	e_dir_usage
n@akamai.com	2.0 2023-0	5-17 14:06:34.166	2023-05-17	5981	15293997		A	/@akamai.com	@akamai.com	_common_data	information_schema	check_constraints	
n@akamai.com	2.0 2023-0	5-17 14:06:17.780	2023-05-17	5981	15293997		A	hushau@akamai.com	data analanarina	eri common data	information cohoma	chack constraints	
n@akamai.com	2.0 2023-0	5-17 14:05:34.756	2023-05-17	5981	15293997		A			1 2 3 4	5 8 >		
© 22 minutes ago								O an hour ago					
Downloads Activitys								Created SQL endpoints					
created_at	workspace_id	source_ip_ad	dress email		action_n	ame note	bookFulli	created_at	created_on	source_ip_address	user_identity		service_nam
2023-06-14 20:22:30.20	7 41223079800004	16	7 1	@akamai.com	workspa	ceExport /Use	rs/jcaruge	2023-06-06 08:26:22.973	2023-06-06	-	{"email":"	1","subjectName":null}	databrickssql
2023-06-14 20:22:30.10	5 412230798 4	16		@akamai.com	workspa	ceExport /Use	ers/jcaruge	2023-06-06 08:26:15.200	2023-06-06	Parate Charles	{"email":"	1°,"subjectName":null}	databrickssql
2023-06-14 20:22:30.02	6 412230798 4	16		@akamai.com	workspa	ceExport /Use	ers/jcaruge	2023-06-15 11:13:21.093	2023-06-15		{"email":"	"subjectName":null}	databrickssql
2023-06-14 20:22:29.75	7 412230798 4	16	r	@akamai.com	workspa	ceExport /Use	rs/jcaruge	2023-06-15 11:14:22.577	2023-06-15		{"email":"	"subjectName":null}	databrickssql
2023-06-14 20:22:29.63	5 412230798 4	16	2	@akamai.com	workspa	ceExport /Use	ers/jcaruge	2023-05-22 06:14:16.094	2023-05-22	1 4	{"email":"	1","subjectName":null}	databrickssql
2023-06-14 20:22:29.53	8 412230798 4	16		@akamai.com	workspa	ceExport /Use	ers/jcaruge	2023-05-22 06:16:25.741	2023-05-22	1 4	{"email":"	1","subjectName":null}	databrickssql
2023-06-14 20:22:29.44	0 412230798 4	16		@akamai.com	workspa	ceExport /Use	ers/icaruqe	2023-05-24 11:25:33.385	2023-05-24	ε	{"email":"	'subjectName":null}	databrickssql
								2023-05-28 08:00:12.958	2023-05-28	10.000	{"email":"	'subjectName":null}	databrickssql
		1 2 3	4 5 40 >					2022 05 20 00:00:42 676	2022 05 20		(1-molt.)	in dianthian star ID	databelahasal

Share Schedule v

O Refresh

### High Level Governance Dashboards

#### **User Privileges Report**

User Privilo	eges Report + Add ta	g								
ir	v									
vot table - Res	ource Privilege - User Privile	ge								
esource_type	resource_name	privilege_type	ALL_PRIVILEGES	CREATE_EXTERNAL_LOCATION	CREATE_EXTERNAL_TABLE	READ_FILES	SELECT	USE_CATALOG	USE_SCHEMA	WRITE_FILES
			ALL_PRIVILEGES							
atalog	( ) a		ALL_PRIVILEGES							
							SELECT	USE_CATALOG	USE_SCHEMA	
					CREATE_EXTERNAL_TABLE	READ_FILES				WRITE_FILES
			ALL_PRIVILEGES		CREATE_EXTERNAL_TABLE	READ_FILES				WRITE_FILES
			ALL_PRIVILEGES							
cternal Location			ALL_PRIVILEGES							
					CREATE_EXTERNAL_TABLE	READ_FILES				WRITE_FILES
					CREATE_EXTERNAL_TABLE	READ_FILES				WRITE_FILES
			ALL_PRIVILEGES							
			ALL_PRIVILEGES							
		ALL_PRIVILEGES								
		ALL_PRIVILEGES								
	A REAL PROPERTY.		ALL_PRIVILEGES							
chema	-		ALL_PRIVILEGES							
	-		ALL_PRIVILEGES							
							SELECT		USE_SCHEMA	
							SELECT		USE_SCHEMA	

41

Akamai

### **Unity Catalog Migration Recipe**

### Putting the pieces together



### Akamai's Unity Catalog Migration Journey

### Foundational building blocks

#### **One Time Setup**

- Identity federation
- Create account groups
- Create account level dashboards

### Akamai's Unity Catalog Migration Journey

### Foundational building blocks

#### **One Time Setup**

- Identity federation
- Create account groups
- Create account level dashboards

#### **Per Region**

- Create a metastore
- Assign metastore to Workspaces
- Create common-data catalog
- Create metastore level utilities & dashboards

### **Akamai's Unity Catalog Migration Journey**

### Foundational building blocks

#### **One Time Setup**

- Identity federation
- Create account groups
- Create account level dashboards

#### **Per Region**

- Create a metastore
- Assign to Workspaces
- Create common-data
   catalog
- Create metastore level utilities & dashboards

#### **Per Product**

- Create catalog, schemas (environments) & tables
- Set product owners
- Migrate legacy metadata & managed tables
- Migrate clusters, jobs & notebooks
- Revoke mounts & Hive\_metastore

### Lessons Learned From the Journey

(So far...)

- Invest in requirements collection early.
- Use automation for efficiency
- Stay updated on new features.
  - Some noteworthy features:
    - Serverless capabilities
    - Audit logs
- Based on our experience, UC performs optimally with Delta Lake format

### What You Learned Today?



### What You Learned Today?



Unity Catalog is essential for Databricks users

- Central access management
- Share data between workspaces easily
- Enhanced security



Starting your journey towards using Unity Catalog

- Collect requirements.
- Derive design decisions based on requirements
- Follow the migration journey slide and do adjustments for your use case

48 @giladasulin, @pulkitchadha

### What You Learned Today?



Unity Catalog is essential for Databricks users

- Central access management
- Share data between workspaces easily
- Enhanced security



Starting your journey towards using Unity Catalog

- Collect requirements.
- Derive design decisions based on requirements
- Follow the migration journey slide and do adjustments for your use case



#### Building utilities using Unity Catalog's metadata

- Catalog as PDF
- State of the Lakehouse Dashboard
- Users' Activity Monitor
- User Privileges Report



### Thank You! Your feedback is important! Feel free to reach out 🙂



