

Healthcare Data with Unity Catalog: Providence's Journey

Anna Erickson, Satish Marriselli, Janet Vickers, Lawrence Yapp


Providence St Joseph Health



122K
CAREGIVERS



38K
NURSES



34K
PHYSICIANS



\$2.1B
COMMUNITY
BENEFIT



51
HOSPITALS



1000
CLINICS



29M
TOTAL
PATIENT VISITS



2.6M
COVERED
LIVES



1700+
PUBLIC
RESEARCH
STUDIES



1
HEALTH
PLAN



18
SUPPORTIVE
HOUSING
FACILITIES



HIGH SCHOOL
NURSING
SCHOOLS &
UNIVERSITIES



178
PARTNER SITES



87
COMMUNITY
CONNECT
PARTNERS



Our Problem – The Wild West

- 131 unique workspaces across 10 subscriptions
- No visibility into cost and usage
- Experimentation with no governance, no cluster policies
- Duplicate data as each workspace had its own source data
- Opportunities for optimization
- Opportunities for tighter collaboration across teams



The Team

We are a cross functional group spanning multiple IS teams.

- Healthcare Intelligence
- Cloud Infrastructure
- Cyber Security
- Networking & Firewall
- And Databricks!



Onsite Sessions

- We prepped and planned for multiple in-person sessions that lasted a week each time
- Focus was on collaboration and achieving a specific set of deliverables
- Because we were in one room, we could quickly clear barriers
- “Let the architects, architect”



Team Building, Even After Hours!

While we were onsite, we continued to collaborate and discussed project plans over dinner.

- It may seem insignificant but social gatherings were needed for team forming and building trust and rapport.



Setting up Infrastructure

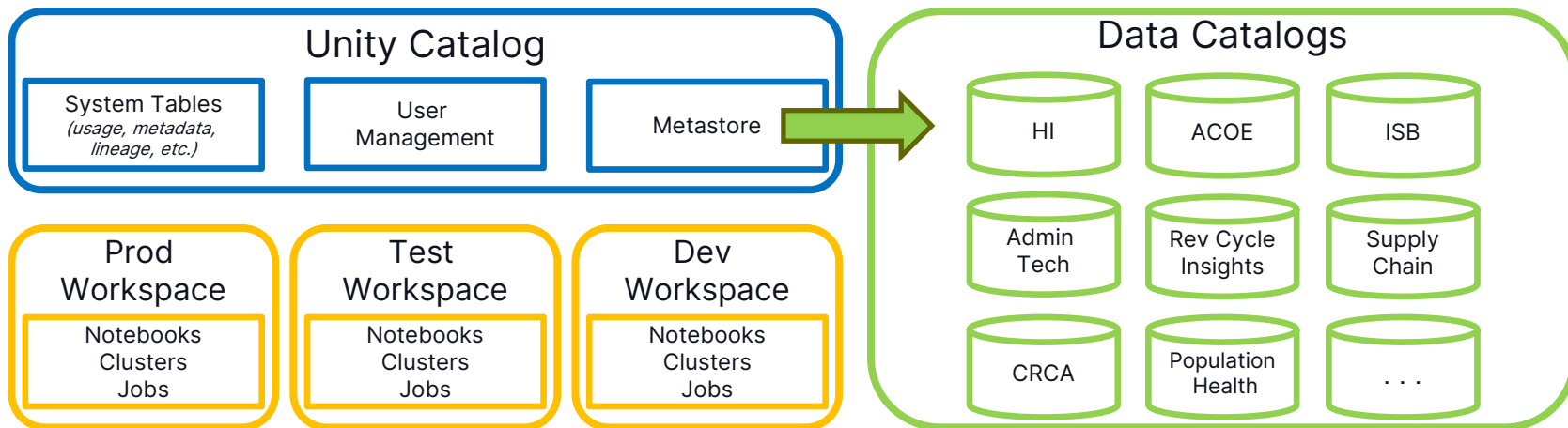


- Networking and firewall configuration
- Involved Cyber Security throughout entire process
- Infrastructure clean up
 - Source data stored in ADLS instead of DBFS, Code Repo and Monitoring
 - Workspace by workspace information gathering
- Deployed single workspace for dev, test, and production in the central subscription

Databricks Unity Catalog Overview

Unity Catalog enabled in a centrally managed subscriptions for the enterprise. This makes it easier to govern, manage users, and permissions.

- Data Catalogs structured similarly to Snowflake Databases
- Metadata layer that interacts with data – schemas, tables, views, etc.
- Data resides in ADLSGen2



Consult with Data Architect to avoid data swamp



Unity Catalog – Benefits and Impact



Why are we consolidating, cleaning up, and migrating to Unity Catalog?

- Reduces the number of workspaces from 131 to 3 (one per environment)
- Centralizing makes it easier to manage and remove duplicate data
- Remove low usage & low value workspaces and clusters
- Easier to secure, govern, and apply policies and manage users/permissions
- Easier to enable CICD and Infrastructure as Code (IaC)
- Easier to monitor and to enable lifecycle management (remove unused clusters, etc.)
- Easier to enable tagging for cost and usage reporting
- Easier to enable new AI features

Migration is not a lift and shift



We began the migration in Jan 2024 and it's a work in progress

- Requires prep work with legacy workspace owners
- Understand existing workflows, data sources and destinations, and how to organize their data in new UC
- Test new jobs and pipelines before shutting down legacy workspaces
- Onboarded 20+ different teams to UC!

HI Migration Counts	Workspace
Starting Count	76
Completed/Removed	48
Percentage Complete	63%

Additional Stats	Workspace
Current Count	33
Consolidated	43
Removed	5
Not Started	1
In Process	32
Remaining	28

NON-HI Migration Counts	Workspace
Starting Count	55
Completed/Removed	33
Percentage Complete	60%

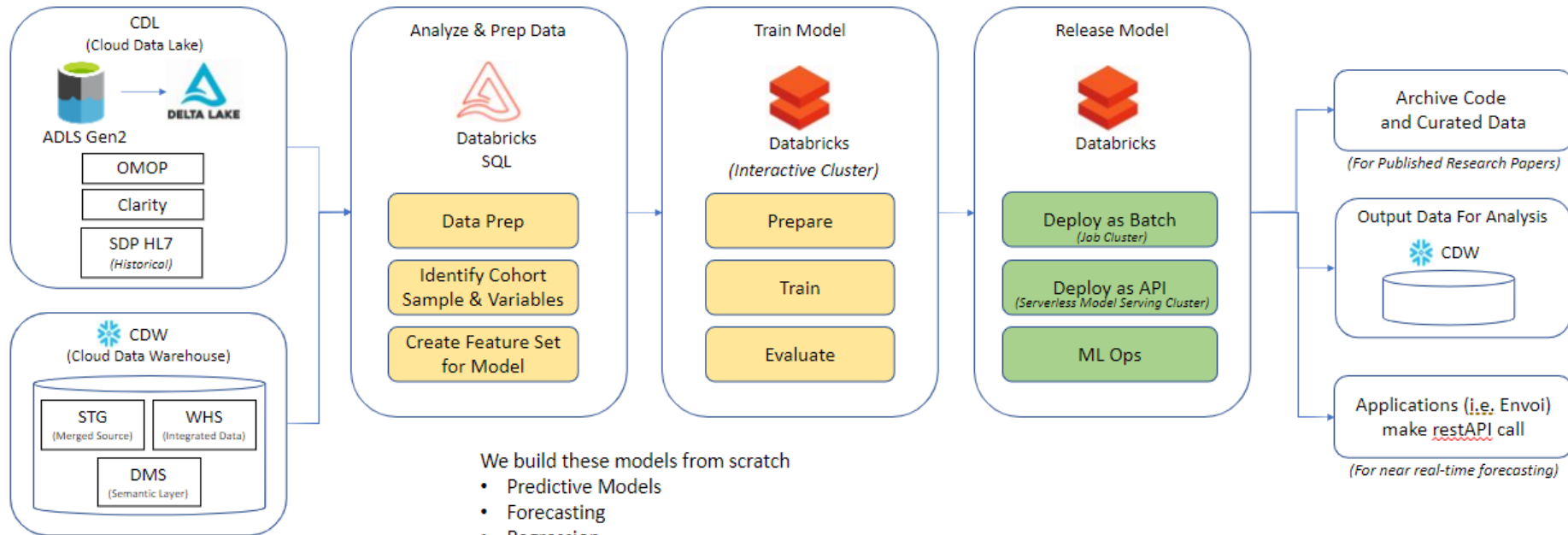
Additional Stats	Workspace
Current Count	39
Consolidated	12
Removed	21
Not Started	15
In Process	3
Remaining	17



Design Patterns - Machine Learning

- Reviewed existing HI Workspaces and created 7 design patterns
- Helped inform cluster policy recommendations and cross collaboration team on use cases

Data Sources

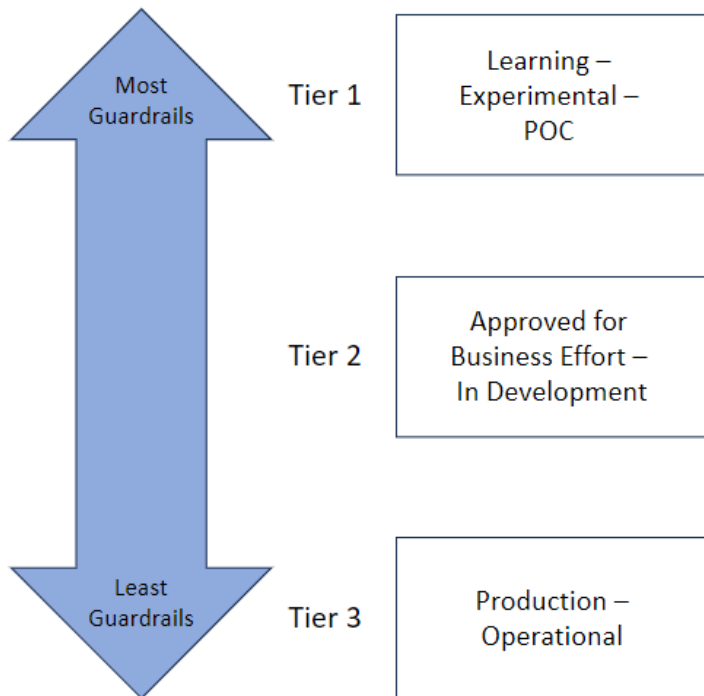


We build these models from scratch

- Predictive Models
- Forecasting
- Regression
- Classification

Enable Guardrails via Policies

- 40 Policies put in place thus far
- Cluster Policies are tied to a users' AAD groups



Cluster Guardrails

Restrict to types of clusters (job, interactive, dlt, etc.)?

Restrict from scheduling interactive clusters

Restrict Interactive Clusters in Prod (only for development)

Set default to small cluster/VM size

Set default to auto-terminate in 30 minutes

Are there certain types of VMs with reservations we want to promote the use of (make default)?

How many machines in the cluster? (min/max range)

Default # nodes

Do we want to restrict dbus used in hour?

Assign tags with default values based on type of processing/pattern

Other Guardrails

Restrict new workspace creation outside of Unity Catalog

Ability to opt-in to new features (deny preview items)

Reinforced Tagging Through Policy

- We used to have one workspace per project
- Now that we're consolidating workspaces, it was important for us to use tags for visibility into telemetry, cost/usage, and lifecycle management.
 - Policies by users' AAD Group, Pattern, and cluster type
 - Automatically assign tags or have defined list of values (no more free text)

Business Customer

Who is this for? Which downstream customer benefits from this processing?

Business Effort

Business friendly name for the project/effort

Project State

Used to indicate where the cluster is in the lifecycle and to manage guardrails.

Support Email

Contact email for operational support

Team Owner

Functional team that manages the cluster

Compute > New compute >

Marripelli, Satish K's Cluster

Databricks runtime version

Runtime: 13.3 LTS (Scala 2.12, Spark 3.4.1)

Use Photon Acceleration

Worker type: Standard_D3_v2 (14 GB Memory, 4 Cores) | Min workers: 2 | Max workers: 8 | Spot instances

Driver type: Same as worker (14 GB Memory, 4 Cores)

Enable autoscaling

Tags

CostCenter	49584503
business_customer	Healthcare Intelligence
business_effort	Prov Chat
project_state	Production/Operational
support_email	HIDataScience@providence.org
team_owner	HI Data Science

Add tags

Key	Value	Add
-----	-------	-----

Cloud Value Optimization Goal

- For Azure Databricks, we saved ~ \$420K in 2023
 - Partnered with Databricks and Microsoft to understand Total Cost of Ownership – DBU, Support, Influenced Costs
 - Reviewed highest cost workspaces
 - Understand value the data processing provides
 - Collaborate with Databricks to identify key data patterns to inform guardrails and provide visibility
 - "Show backs" leading to changes in behavior
 - Identified low value workspaces that should be sunset
 - Focus on optimization
 - Changing timeouts and right-sizing high cost VMs
 - Optimizing Foundational Data (Clarity Ingest, De-Identification, Truveta Trucking Service)
 - Remove low cost/unused workspaces
 - Use data to better inform P3 Reservations and future commit agreements
- Continued focus on optimization in 2024
 - Unity Catalog Migration: focus on rightsizing VMs, optimization, and turning off low value clusters.
 - Interactive vs Job Clusters: policy put in place to restrict scheduling clusters (up to 5X the cost of job clusters). Work on changing behavior to use interactive clusters for development only.

Databricks Total Cost of Ownership

When a Databricks cluster is running, it incurs DBU, Support, and Influenced Costs. Use Azure Portal for big picture view of total cost of ownership.

Central
Managed
Subscription

rg-dbx-enterprise-dev-wus2

Non-Managed Resource Group Includes Databricks DBU + Support Costs

rg-mgmd-dbx-enterprise-dev-wus2

Managed resource group includes Databricks Influenced Costs

Resource Group	Databricks Cost Type	Service Name	Amortized Cost Amt*	% TCO
rg-dbx-enterprise-dev-wus2	DBU	Azure Databricks	126,057	63.03%
	Support	Virtual Network	223	0.11%
		Azure DNS	10	0.00%
rg-mgmd-dbx-enterprise-dev-wus2	Influenced	Virtual Machines	52,375	26.19%
		Virtual Network	13,205	6.60%
		Storage	7,901	3.95%
		Bandwidth	178	0.09%
		Microsoft Defender for Cloud	50	0.02%
Grand Total			200,000	100%

Note: Serverless clusters will not show Influenced Costs in Azure Portal. Expect DBU spend to increase to reflect these costs.

* Synthetic data



Cost Optimization - Rightsizing VMs

Pages << File Export Share Chat in Teams Get insights Subscribe to report ...

DBX Amortized Cost

Prod VMs Utilization

Dev VMs Utilization

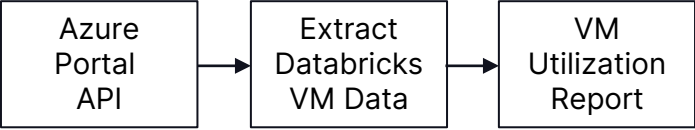
Clusters Amortized Cost

Other Services Amortize...

Databricks VMs CPU & Memory Usage

VM Avg CPU % and Available Memory Details						
ResourceGroup	VM Name	Size	DateTime	Avg CPU %	Available Memory (GB)	
RG-DBX-NOSHOW-MANAGED	90272627da004e08935847d7245d386	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/19/2024 3:02:06 AM	0.65	361.54	
RG-DBX-NOSHOW-MANAGED	a9ebf2a053d24e4f93361b67b56f406f	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/19/2024 3:02:08 AM	0.54	347.23	
RG-DBX-NOSHOW-MANAGED	9d68329c50184039bc2affb4e807d6d0	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/2/2024 5:03:22 PM	0.33	233.11	
RG-DBX-NOSHOW-MANAGED	930b319bca8241749f41fce2fd82ba63	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/2/2024 4:03:51 PM	0.30	213.47	
RG-DBX-NOSHOW-MANAGED	9d68329c50184039bc2affb4e807d6d0	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/2/2024 4:03:53 PM	0.30	213.15	
RG-DBX-NOSHOW-MANAGED	19d5275f8cf04c27af51a9ae9ee661bb	Standard_D64s_v3	5/18/2023 4:18:47 AM	0.59	211.83	
RG-DBX-NOSHOW-MANAGED	a1d92ac7cc144c6f8a67e4cc2ca3f758	Standard_D64s_v3	5/18/2023 4:18:53 AM	0.60	210.38	
RG-DBX-NOSHOW-MANAGED	99615a7e5dd14203a008a45eac67e5e8	Standard_D64s_v3	5/18/2023 4:18:52 AM	0.56	209.10	
RG-DBX-NOSHOW-MANAGED	4e0fda25d3ac442291c7a364681c2b8a	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	2/5/2024 12:09:45 PM	0.18	179.85	
RG-DBX-NOSHOW-MANAGED	a9bfa437db024105b3a7c70641878716	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	3/16/2024 10:00:59 PM	0.12	172.34	
RG-DBX-NOSHOW-MANAGED	90272627da004e08935847d7245d386	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/19/2024 2:02:17 AM	0.62	171.23	
RG-DBX-NOSHOW-MANAGED	a9ebf2a053d24e4f93361b67b56f406f	Standard_D64s_v3 (64 vCPUs, 256 GiB memory)	1/19/2024 2:02:21 AM	0.52	164.94	
RG-DBX-NOSHOW-MANAGED	90272627da004e08935847d7245d386				161.53	
RG-DBX-NOSHOW-MANAGED	90272627da004e08935847d7245d386			0.12	172.34	
RG-DBX-NOSHOW-MANAGED	a9ebf2a053d24e4f93361b67b56f406f				155.34	
RG-DBX-NOSHOW-MANAGED	a9ebf2a053d24e4f93361b67b56f406f				155.34	
Total					73,584.93	

Extract data to find oversized and underutilized Databricks Virtual Machines

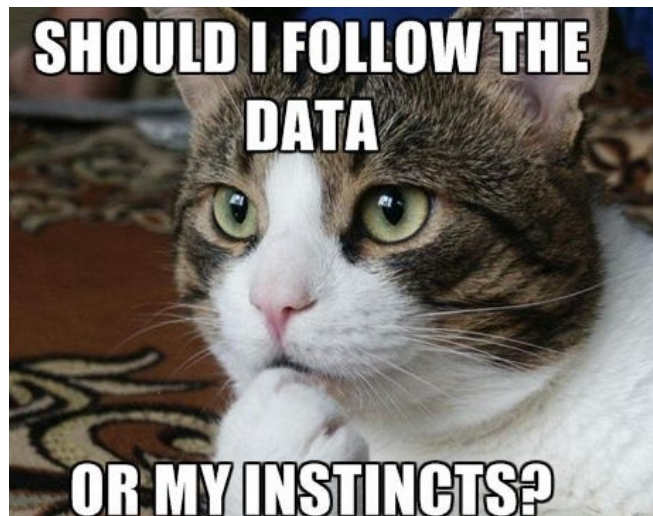


Low Average CPU % Utilized and High Available Memory (GB) indicated we need to follow up with the project teams to right size the clusters/VMs.

Unity Catalog - System Tables

Allows for reporting of both Unity Catalog and non-Unity Catalog workspaces

- Identified clusters with time out of 120+ minutes
- Jobs usage (the exact DBUs that each jobs cluster consumes)
- DBUs by cluster type/SKU
- Audit logs – tracks user activities



What about serverless? What's Next?

We are focused on serverless and other Databricks features that advance our AI work.

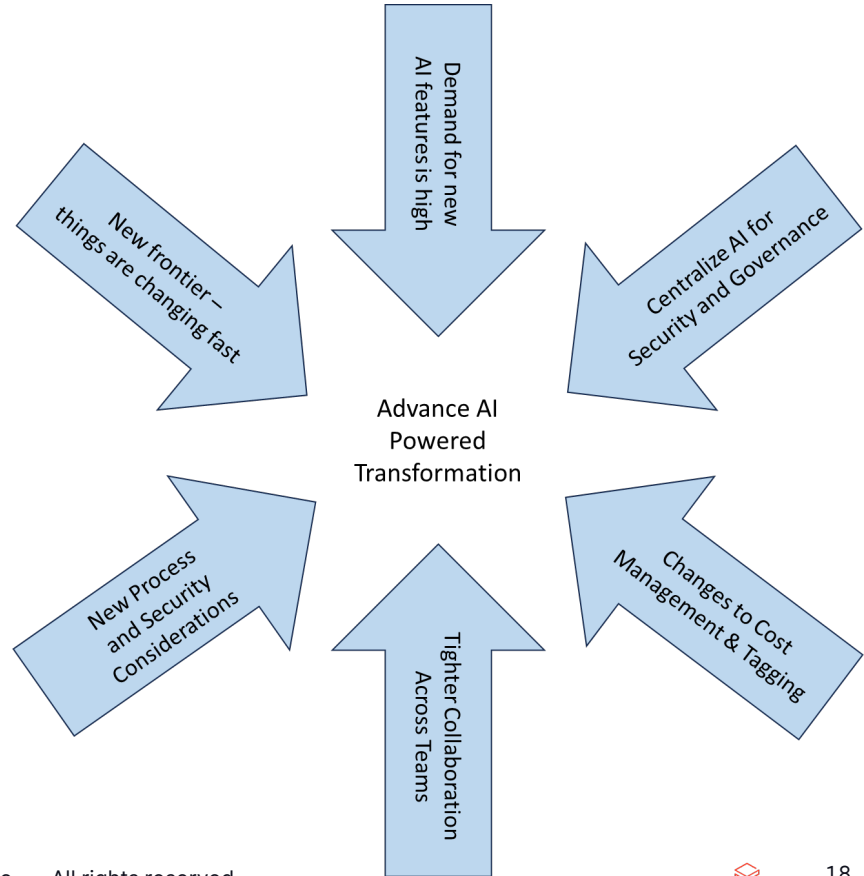
Enabled Serverless Model Serving

- Initial use case is to deploy one of our machine learning models as a restAPI so it can be called by downstream applications

Emerging AI Forces

Healthcare is a highly regulated environment.
Tension between:

Being Innovative  Being Safe/Secure



Enabling AI Features

- Partner with Databricks to share our AI use cases and learn about upcoming features
- We are focused on infrastructure to enable Azure Open AI Models
- Keep Databricks Feature Road Map up to date and plan quarterly for next set of features to enable



Questions?

Thank you for attending our presentation!

