



Episode 53

Build the next-generation Iceberg lakehouse with Dremio and NetApp



July 9th, 2024 | 8AM PDT | 11AM EST | 4PM BDT



Vishnu Vardhan

Director of Product Management StorageGRID NetApp



Alex Merced Senior Technical Evangelist Dremio

Dremio: Expert and Leader in the Apache Iceberg Lakehouse



- One of the first platforms to embrace Apache lceberg
- Created some of the earliest online Apache Iceberg education materials
- Creators of O'reilly's "Apache Iceberg: The Definitive Guide"
- Creators of the first open-source Apache Iceberg Catalog (Nessie), and only catalog with catalog level Git-for-data features.



An Apache Iceberg Crash Course 10 Sessions July 11 - October 29



Alex Merced Sr. Evangelist Dremio







Agenda

- Market News
- Why Iceberg Lakehouses
- Dremio & NetApp Joint Offering
- NetApp StorageGRID & Dremio Deep Dive
- Joint Solution Benefits & Use Cases
- Case Studies



Data Lifecycle Remains Complex, Brittle, and Expensive





Dremio & NetApp Iceberg Data Lakehouses Help!



Reduced Cost, no ETL to a data warehouse, reduced extracts, no copies to manage Simplifies data pipelines, improves "uptime" for data accelerates and time-to-insight



Databricks purchases Tabular

Deals

Databricks to buy data management firm Tabular for over \$1 bln

By Reuters

NEWS

June 4, 2024 11:39 AM PDT · Updated 17 days ago

June 4 (Reuters) - Databricks said on Tuesday it would buy data-management startup Tabular for more than \$1 billion, as the privately held analytics platform looks to attract customers by helping them develop custom artificial intelligence (AI) ap

The announcement comes at a time clients to use open-source AI mode <u>(SNOW.N)</u> and Cloudera, among



Databricks Inc. today agreed to acquire Tabular Technologies Inc., developer of a universal storage platform based on the Apache Iceberg standard.

Aa

The move signals stepped-up efforts by Databricks to bridge the compatibility gap between its Delta Lake storage format and Iceberg. Terms weren't announced, but Databricks Chief Executive Ali Ghodsi (pictured) told CNBC the price tag was more than \$1 billion. Snowflake Inc. and Confluent were also reportedly in on the bidding.

Databricks \$1B-plus Tajahalar was founded by three former Netflix Inc. employees who co-created Iceberg at that company. support

The lakehouse specialist's latest purchase adds support for Apache Iceberg to its existing support for Delta Lake and is also a direct confrontation of rival Snowflake.

Databricks continued its recent buying spree with the acquisition of Tabular, a move that adds support for Apache lceberg storage to Databricks' existing support for Delta Lake storage.

Databricks did not specify what it paid for Tabular but confirmed that the cost was more than \$1 billion.



HPE Discover

HPE GreenLake

Home \rightarrow Al/ML \rightarrow Databricks buys Tabular to win the Iceberg war

AI/ML Data Management

Databricks buys Tabular to win the Iceberg war

By Chris Mellor - June 5, 2024





APP DEV CLOUD GEN AI MACHINE LEARNING ANALYTICS NEWSLETTERS RESOURCES

Databricks to acquire storage platform maker Tabular

While Snowflake is talking up its use of Iceberg to promote interoperability, Databricks is buying Tabular, the tool built on Iceberg's table format by Iceberg's creators.





Snowflake Releases Polaris Iceberg Catalog

Snowflake Releases Polaris Catalog: Transforming Data Interoperability with Open Source Apache Iceberg Integration

By Asif Razzaq - June 4, 2024



Snowflake has unveiled the **Polaris Catalog**, an open-source catalog for Apache Iceberg that enhances data interoperability across various engines and cloud services. This launch signifies Snowflake's commitment to providing enterprises more control, flexibility, and security for their data management needs.

The data industry has increasingly embraced open-source file and table formats for their potential to improve interoperability. This capability allows multiple technologies to operat over a single copy of data, reducing complexity, costs, and risks associated with vendor

DATA

Snowflake, Databricks and the Fight for Apache Iceberg Tables

The market for data lakes and data lakehouses is clearly being disrupted by open source software, given recent news from Databricks and Snowflake. Jun 10th, 2024 9:46am by Joba Jackson

CLOUD SERVICES / DATA

Snowflake Polaris Aims for Multiquery Engine Interoperability

An open source catalog for Apache Iceberg, Snowflake Polaris provides a way for multiple query engines to write Apache Iceberg tables.

Snowflake unveils Polaris, a vendor-neutral open catalog implementation for Apache Iceberg





vflake catalog supports cross-engine acces

BY PAUL GILLIN



June 3, 2024

Snowflake Embraces Open Data with Polaris Catalog $_{\rm Alex\ Woodie}$



cello/Shutterstock)

On the first day of its Data Cloud Summit toda Snowflake unveiled Polaris, a new data catalog data stored in the Apache Iceberg format. In addition to contributing Polaris to the open so community, the catalog also enables Snowflake customers to use open compute engines with Iceberg-based Snowflake data, including Apach Spark, Apache Flink, Presto, Trino, and Dremio.

TECHNOLOGIES -

The launch of Polaris represents a significant embrace of open source and open data on the of <u>Snowflake</u>, which grew its business predomin through a closed data stack, including propriet table format and a proprietory SOL processing

Snowflake adopts open source strategy to grab data catalog mind share

With its plan to make its Polaris data catalog open source, Snowflake hopes the new offering will be seen as vendor-neutral, boosting its attractiveness when compared to Databricks' Unity Catalog.

() 🖸 🗇 😋 🕞



y Anirban Ghoshal nior Writer, InfoWorld | JUN 3, 2024 6:00 AM PDT

Jun 3rd, 2024 2:05pm by Joab Jackson

Why Iceberg? A Diverse Developer Community



% Attributable Contributions to Apache Iceberg by Company





% Contribution to Delta Lake by Company





Iceberg for Your Next Generation Lakehouse

Queries are faster (even without perfect users)

- Automatic up-to-date statistics
- Hidden partitioning reduce full table scans
- Optimal data layout (file sizes, prefixes, etc.)

Changes are easier

- No need to refresh metadata after updating a table
- Atomic updates in one command (INSERT/UPDATE/DELETE)
- Raw files (CSV/JSON) to queryable and transactional table in one command (COPY INTO)
- Schema/partition evolution with one command

Automatic data optimization

Automatic garbage collection

Less work for

- No metadata refresh jobs to trigger and manage
- No manual work to insert/update data or evolve schema/partitions
- Recovering from mistakes and disasters is trivial
- Fewer tickets!!

Lower compute and storage costs

- Time travel reduces data copies
- Less compute per query (see "queries are faster")
- No need to copy data into a data warehouse

GNARLY Data_Waves

- + No vendor lock-in
- + Broadest ecosystem of any table format

NetApp StorageGRID

Industry leading Object Storage solution

- Single Global Active/Active Namespace
- Powerful Policy Engines
 - Automatically place data to the selected media, data center, and/or protection scheme
- Flexible & Simple to Deploy
 - Mix & match software-defined and appliances
- Unmatched Durability
 - Achieve 15 9's of durability leveraging layered erasure coding
- Scalable Architecture
 - Supports 300 billion objects and 640+PBs in single namespace
- Enabling the Hybrid Cloud
 - Cloud integration with AWS workflows and services
 - Tier data to Microsoft Azure, AWS S3/Glacier & Google Cloud





Automate Lakehouse data management at scale with ILM policies

Policies determine

- Where data lives
 - Place objects in specific data centers or specific media (disk, tape, cloud)
- How long data is stored
 - Have different lifecycles for objects
- How secure the data is stored
 - Create more or fewer replicas of objects or erasure encode objects

Apply policies on existing objects

 StorageGRID[®] automatically brings objects into policy alignment with a fast background engine

Name				In Active Po	blicy	In Proposed Polic
Make 2 Copies				~	•	
German Data Policy						
						4
German Data Polic	y					Version 1.0
Description:	Data	needs to sta	ay in Gern	nan Datacenters		
Tenant Account:	9435	2640785552	2504147			
For Object Type:	S3/S	wift				
Reference Time: Filtering Criteria:	Inges	st Time				
Matches all of the fe	ollowing me	tadata:				
location		equals	1	germany		
Retention Diagram:						
Trigger			Da	/ 0		
	Frankfurt DC		9			
	Berlin DC		9			
Duration				Forever		



Simple and feature rich software-defined object storage





Dremio: The Unified Lakehouse Platform for Self-Service Analytics & Al





RESENTED BY **Odremio**

Unified Analytics

	BI tools, data science notebooks, SQL editors							
	Ļţ	ODBC JDBC REST Arrow Flight	↓ ↑					
Gert	Self-Service Analytics	 Shareable, governed data Views GenAlText-to-SQL Comprehensive BI integration Embedded SQL runner 	 Governance & Security Role-Based Access (RBAC) Fine-grained access control 					
UNIFIED ANALY TICS	Semantic Layer	 Business-focused virtual data sets and marts Seamless search Intuitive, user-generated project Wikis 	 Auditing & Query history Data Lineage Identify Management & SSO integration 					
Leberg, Delta Lake Parquet, ORC, JSON, CSV Readers ARP Connectors								



SQL Query Engine

BI tools, data science notebooks, SQL editors

	↓ ↑	ODBC JDBC REST Arrow Flight	↓↑
Gart	Price Performance Reflections	 Auto-scaling/elastic engines Multi-engine architecture Workload management, query routing Reflections Acceleration via optimized relational cache 	Multi-Cloud & Hybrid • On-premise • In the cloud
SQL QUERY ENGINE	Query Acceleration	 Automated Reflections Recommender Columnar Cloud Cache (C3) Query federation 	 Multi-cloud Hybrid architecture
	Federation	Connector Ecosystem	

Iceberg, Delta Lake | Parquet, ORC, JSON, CSV Readers | ARP Connectors



Lakehouse Management





Dremio Hybrid Iceberg Lakehouse



- Enterprise Hybrid Cloud Iceberg Catalog
- Leverage on-prem and cloud based storage
- Enable self service across all data sources



Why Dremio + NetApp



Iceberg Ingestion Options with Dremio



Dremio + NetApp StorageGRID Iceberg Lakehouse



Dremio & StorageGRID Automate Data Management at Scale

Isolation, version control, governance & table optimization

Dremio Manages:

- User Access
 - Fine-grained privileges to control access to the data at the table, column and row level
- Version Control
 - Recover from any mistake by instantly undoing accidental data or metadata changes
- Table Optimization
 - Automatically compacts small files and group similar rows together
- Isolation
 - Enables experimentation with data without impacting other users

🛄 Drem	io Alliances / 🔌 Sona	r / 👔 Alliar	nces	✓ / 🖂 S	QL Runn	er						
Data	Scripts	$\leftarrow \mid$	(May 28, 2024,	11:46:35	AM X May 28, 202	24, 11:47:05 AM	× Jun 3, 2024, 12:43	3:12 PM × +		🛅 nyc_trips ϟ	\mapsto
All	Starred (0)	Name ↑		🕑 Run	٢	Engine: 🗞 Automati	c ~		× (2)	Save as View 🗸 🗸	~ Overview	
> 🎧	@mark.shainman@dre	mio.com							Context: (None selected)	fx 🔅 🖸 🗐	NetApp."dremio-b	oucket"."nyc_trips"
	Airbyte			1 SELECT	* FROM	NetApp."dremio-bu	cket"."nyc_tri	ps"			No Label 🖉	
>	Alliances Alliances S3										Jobs (last 30 days	s) 2
> 💽	arctic 🕼 main >										Descendants	0
	co-sell aws										Created	04/17/2024, 06:28:46
> 🔲	Curated Data										Owner	RO rommel.garcia@
	dbt Datasets										Last updated	06/03/2024, 12:40:25
	Demo Data										Launch BI tool	* 4
> 8	Demo Data 2										Go to	N 29
> 🔲	dev										v Columns (20)	0
> 🖯	Fivetran Brown Bag										• Columns (20)	4
> 🖯	Fivetran Iceberg										abc vendor_id	
> 😫	Fivetran_Dremio_Glue										🛱 pickup_date	
> 😫	NetApp										🛱 pickup datetin	ne
> 😫	On-Prem Object Stora	ge										
> 🔲	Product					Filter Columns	0 Columns				H dropoff_date	
> 🛛	Production Data										🗄 dropoff_datetir	me
> 🔲	Raw Data											
> 🔲	Sales						Dura a O				∽ Wiki	
) A	SampleDB						Run a Query	to Get Started				



Benchmark Results: NetApp StorageGRID & Dremio win!

	80% range read 2KB 2MB from 32MB obje 50 – 100 requests/s	B –73% range read belowects,100KB from 32MB objectsec1000 – 1400 requests/se	90% 1M byte range read from 256MB objects, 2000 – 2300 requests/sec	
S3 Requests	Hive	Spark + Delta Lake	Dremio 4,414,227	
GET (range read)	1,117,184	2,074,610		
List Objects	312,053	24,158	240	
HEAD (non-existent object)	156,027	12,103	192	
HEAD (existent object)	982,126	922,732	1,845	
Total Requests	2,567,390	3,033,603	4,416,504	
99 Queries: Total Time (Minutes)	1084*	55	47	
		* Hiv	ve unable to complete query #72 in this test	

The unable to complete query $\pi I \ge 111$ this test

The NetApp & Dremio Difference





Dremio & NetApp: Hadoop Modernization

- Sub-second query performance & 10x better price/performance
- Governed self-service analytics
- Unified view & access across all of an organization's data
- Ability to scale compute and storage separately
- Improved overall data management
- Reduced complexity

Migrate from a legacy Hadoop data lake to a modern NetApp StorageGRID & Dremio lakehouse environment!



Superior business insights & faster time to value!



Use Case: Hadoop Modernization for NetApp's Own Lakehouse!

Product 360: 95% faster time-to-insight while simplifying proactive customer care

Solution Overview

 Active IQ is NetApp's Digital Advisory solution for predictive maintenance and optimization



Pain Points

- Poor performance with legacy Cloudera/Hadoop technology
- Storage and compute tightly-coupled, leading to increased costs
- Scaling performance to meet workload demands was challenging: the average Hive query took more than 45 mins
- Dealing with costs of having to copy and move data between data sources

Why NetApp StorageGRID & Dremio

- Ability to guery data directly in the datalake, reducing need to copy or move data
- Reuse existing investment of compute & storage, minimize changes to existing data pipelines
- Queries are 95% faster with Dremio & StorageGRID: 45 min down to 2 minutes!
- 30+% analytics TCO savings in Year One!

Results

GNARLY Data Waves

PRESENTED BY CORNIG

Empowered data consumers with self-service layer for BI team to access data

Proactively manage customer care, optimize and identify problems before they happen.

Data

Science

ETL/Data Ingestion

Spark

jupyter SQL



customer churn and higher product

North State

3 PB

Active IQ Data Lake

Dremio on K8S

16 Executor nodes on k8s cluster

8900+ Tables

Apps

SGRID

StorageGRID & Dremio Lakehouse Solution

垫

Unified Analytics & SQL Query Engine

dremio

E SERIES

Data Lakehouse

Dashboards

ONTAP

Use Case: Hadoop Modernization: Global 100 Financial Institution

- Existing Large, Legacy On-Prem Hadoop Data Lake
 - No plans to move it to the cloud
 - Concerns with regulations, security & PII
- Uses Dremio for their Financial Data Landing Zone
 - "Single source of truth" for all of their corporate treasury, liquidity, etc.
 - Able to join and query data from NetApp StorageGRID and existing Hadoop for high-performance BI & reporting
 - With Dremio semantic layer, they're able to provide a self-service analytics layer and make cross-functional data easily-discoverable
 - Dremio makes the audit trail of data easier for mandatory reporting to the FTC
- Long-Term NetApp Plans
 - Reduce Cloudera Data Platform footprint and move off of DAS
 - As data moves off of DAS, it will go to NetApp StorageGRID
 - Majority of Treasury group's new data moving to NetApp StorageGRID
- Long-Term Dremio Plans
 - Adding cloud data sources & relational database data sources with their on-prem data
 - Continue to expand existing on-prem footprint







Use Case: Warehouse to Lakehouse

Fortune 10 Customer: 75% TCO Savings, \$3M Savings In Just One Dept



The NetApp & Dremio Difference

Best-in-Class Best **Price-Performance** TCO **Simplified Data** Fastest **Engineering & Query Performance** Management



Get More Great Information



NetApp®



NetApp Case Study



Dremio & Storage Grid Solution





Five Factors to Consider When Migrating from Hadoop to the Data Lakehouse

Thank You

