

EPISODE XX

What's New in Dremio: Improved Automation, Performance + Catalog for Iceberg Lakehouses

Today's Agenda

Mark Shainman & Maeve Donovan



- Dremio Release Focus
- What's New in Dremio
 - Improved Performance & Automation
 - Improved Data Ingestion
 - Accelerated Cross-Database Access
 Control and Workload Management
 - Integrated Observability
 - Google BigQuery Connector
 - Graviton Support
- Iceberg Catalog

We're making Iceberg Lakehouses more manageable and higher performing than ever before



Decrease TCO and Improve Time to Value

Apache Iceberg

An Open Table Format for Enterprise Data Lakes

High-performance queries - Purpose-built for high performance queries on massive datasets.

Data warehouse functionality on the data lake - ACID transactions, time travel, and schema evolution enable more data warehouse workloads directly on data lake storage.

Easy data operations - Reduce overhead costs with table optimization, garbage cleanup, and more.

The Largest Open Source Community

More individual companies with contributions than any other open table format

More OSS integrations than any other open table format.

Enterprise Companies Using Iceberg NETFLIX Expedia stripe **airbnb** (ii) twilio Linked in A Adobe Tencent **Commercial Support for Iceberg** aws **dremio k**snowflake Google Cloud CLOUDERA ₩ Starburst



Improved Performance and Automation



Reflections Eliminate the Need for BI Extracts & Imports





Reflection Recommendations

- Simplifies the workload of data engineers and enhance the performance of frequently executed queries on views by recommending raw Reflections with high ROI.
- Recommendations are generated daily, and are based on the query history of the preceding seven days.
- ROI is calculated based on the number of prior jobs that could be accelerated and the expected average improvement in performance.

Reflection Recommendations

-	Dremio / Sonar / Dremio Cloud / Project Settings										
▦	Dremio C	Reflection recommendations		\sim							
Þ	(i) Ge	Name	Туре	Dataset	Accelerated job c	Queries speedup	Queries time savi	Manage	Columns		
	Prc	AutoRef_edition_current_raw	📑 Raw	edition_current	113	10.936x	53m 33.5s		Record		
謳		AutoRef_jobs_all_raw	🛃 Raw	jobs_all	23	1.134x	6.6s	13:00:31	4472		
ŝ	2 Eng							13:29:11	84,724,		
	40							13:11:3:	L 9,5		
								13:29:33	37,414,		
								13:14:33	L 3,:		
								13:13:43	3,		
Ē											
_											
?											
SH								/			
То а	ccess th	ne list of recommenda	ations fo	or your project and	d add new Reflectio	ons:					
	Clic	k the Project Setting	s icon i	n the side navigat	ion har		To add a	Reflection to y	our 🧹		
	One		5 100111	in the slot havigat	1011 July 202		Reflectior	ns list, click at	the		
	· On	your Project Settings	page, o	click Reflections .			end of the	e row. 🔶			
_	· To a	access your list of rec	commer	ndations Click 📩	1_						
				+	V +⊳		\triangleright				

Live Reflections

- Simplifies the process of Reflection management and reduce administrative burden for data engineers.
- Reflections based on Iceberg tables are now automatically updated when their underlying tables are updated using the Dremio engine or other external engines.
- Live Reflections on Iceberg table can be enabled at source or at table level.

tings for NYC-taxi-t	rips									
Refresh	Refresh Policy How often reflections will not be refreshed and how long data can be served before expiration. Failed reflections will not be refreshed. Refresh now Refresh now									
	Never refresh Never refresh Refresh very 1 Hour(s) > Starten will not every Monday, sta SM. The next job is scheduled at Sep 23, 2023, and the scheduled at Sep 23, 2023,									
	Cancel Save									

Folder Se

Privilege

Reflection Scores

- Makes it easier for platform engineers to identify and take action on Reflections whose value is low or has dropped.
 Scores are based on the
- Scores are based on the usage and acceleration speed up in the last 7 days.

	Test Project	4 R	eflections						.∰ ¢3
	(i) General information	Q	Search reflection name/id	id/dataset Acceleration Status: All 🗸		Refresh Status: All 🗸 🗸	Mode: All 🗸	0 Manage Columns	
	Project Storage		Name	Туре	Mode	Dataset	Reflection Score	Current Footprint	Last Refr
	📽 Engines	0	Raw_Northwest_1	Raw	Manual	Northwest	11 Poor	8.12 KB	01/19/2
=	III BI applications	43	Agg_Northwest_2	😭 Aggr	regation Autonomous	Northwest	71 Great	8.12 KB	01/19/2
_	B Monitor	0	Agg_Northwest3_1	Aggr	regation Manual	Northwest3	9 Poor	8.12 KB	01/19/2
<u>Ş</u> 3	47 Reflections	0	Agg_Northwest4_1	😭 Aggr	regation Autonomous	Northwest4	78 Great	8.12 KB	01/19/2
	C Engine routing	0	my reflection	Raw	Manual	Northwest5	40 Fair	8.12 KB	01/19/2
	Preference	0	Agg_Northwest6_1	😭 Aggr	regation Autonomous	Northwest6	82 Great	8.12 KB	01/19/2
		0	Raw_Northwest7_1	Raw	Manual	Northwest7	3 Poor	8.12 KB	01/19/2
		43	Raw_Northwest8_1	📑 Raw	Autonomous	Northwest8	72 Great	8.12 KB	01/19/2
		0	Agg_Northwest9_1	😭 Aggr	regation Autonomous	Northwest9	93 Great	8.12 KB	01/19/2
		0	Agg_NYC_1	😭 Aggr	regation Autonomous	NYC	93 Great	8.12 KB	01/19/2
		0	Agg_Sales Data	😭 Aggr	regation Autonomous	🕎 Sales Data	93 Great	8.12 KB	01/19/2
		0	Agg_Sales Data_1	😭 Aggr	regation Autonomous	🔢 Sales Data	93 Great	8.12 KB	01/19/2

Seamless Metadata Refresh

- Queries on Iceberg tables, will now do inline metadata refresh to provide users with the most current results at no performance penalty.
- It relieves data and platform engineers from the burden of setting appropriate metadata intervals and/or triggering metadata refresh.
- Reduces compute cost by not performing redundant metadata refresh on tables that are not frequently queried or updated.

Results Cache

- Accelerates the performance of queries that have been previously executed.
- Caches result sets for JDBC, ODBC and Flight type queries.
- Caches up to 20MB result sets and lasts for 24hrs.
- Stores the data in the customers distributed storage (S3 etc.)
- Benchmarking has shown a 28x increase in performance.
- If the cache is leveraged the plan and visual profile will reflect this.

Merge on Read

- New Iceberg V2 table property
- Increases write speed by writing deletes separately.
- Deleted records are tracked via delete files.
- On read, Data files and deleted records are ... merged.



Merge-on-Read: New Delete File created after DELETE operation



Auto Ingestion for Iceberg Tables

Auto Ingest Pipes

- Reduces time and effort to build/operate/maintain ingestion pipelines for the Data Engineer.
- Quicker access to date reducing latency to near real time (~10min from file creation) for the Data Analyst.
- Automatically ingests data from S3 into Iceberg tables.
- Available for both Dremio Cloud and Dremio Software.
- Pipes have an auto deduplication capability and a default deduplication lookback period of 14 days.
- Auto Ingest Pipes leverage the cloud providers event driven architecture for providing immediate processing and ingestion.



ICEBERG

Auto Ingest Pipes





Accelerated Cross-Database Access Control and Workload Management



Oracle DB, Teradata and Microsoft SQL Server Impersonation

- Queries submitted by Dremio can be configured to run under the users own username.
- Allows for existing access control rules to apply to queries coming from Dremio.
- More visibility of who is running what query for Administrators of respective systems.
- Allows for database native workload management capabilities to be leveraged more easily i.e. Teradata's TASM, Oracle's Automatic Workload Management, SQL Server's Resource Governor









Integrated Observability



Administration and Monitoring with Integrated Observability

Out-of-the-box metrics for Resources

- Admins can now see the fluctuation of cpu and memory of engines on Software to understand peak times, need for scaling, etc.
- They also have a view of the top 10 queries that are most memory or cpu intensive.
 This is in addition to the Catalog and Jobs observability that was shipped in 25.0.



Top 10 CPU intensive jobs	Lest 24	hours v	Top 10 memory intensive jobs		Lest 24 hours	~
Job ID	Duration	CPU	Job ID	Duration	Memory	
e228d8e44f88	00:00:31	00:01:30	3bc996da6a88	00:00:09	781.46 MB	
6f17d6f26788	00:00:20	00:00:20	.241aa2184988	00:00:09	781.46 MB	
📀33364f28b988	00:00:13	00:00:11	_f22ec9e34500	00:00:09	781.46 MB	
🕗ba6a78298488	00:00:14	00:00:11	.264d586c0700	00:00:09	781.46 MB	
O769ba48a5488	00:00:03	00:00:10	.03e25f9d6588	00:00:09	781.46 MB	
d02a6938eb88	00:01:39	00:00:10	_661812985788	00:00:09	781.46 MB	
📀fe2dbaaa1808	00:00:11	00:00:09	2dc97c128188	00:00:09	781.46 MB	
Ø76d8157d1500	00:00:11	00:00:09	_ea09c801eb00	00:00:09	781.46 MB	
9a92c18bc888	00:00:11	00:00:09	📀 _b314db196b88	00:00:09	781.46 MB	
22ce3a3d2288	00:00:03	00:00:09	_543f6ec9f988	00:00:08	781.46 MB	



Google BigQuery Connector

Google BigQuery Connector now in Preview

Enabled via a support key

- JDBC based connector for Google BigQuery
- As with other non-lake sources allows for the federation of datasets in a single query
- Where possible it pushes down the query to BigQuery to minimize data pulled across the wire
- JDBC driver must be installed by the customer, like with the Teradata connector





Graviton Support



Graviton Support

- AWS Graviton is a family of ARM based CPU's designed by Amazon
- Delivered in Amazon EC2
- Provide a lower energy use for workloads and a lower cost
- Dremio now supports Graviton instances in Dremio Software
- Dremio customers now has more choices when it comes to the instance type for their workload profile



Apache Polaris (Incubating)



What is Apache Polaris (incubating)?

Apache Polaris (incubating) is a catalog implementation for Apache Iceberg built on the open source Apache Iceberg REST protocol that allows for:

- Cross-engine read and write interoperability
- Centralized access across engines
- Vendor-agnostic flexibility

What this means for customers:

- Maintain one data copy that many engines can query and write to
- Streamline data management with centralized access
- Run anywhere without lock-in

POLARIS CATALOG Snowflake Snowflake Managed Service EC2 EKS Docker Google Cloud Kicrosoft Azure CEC2 GCE GKE Docker

Apache Polaris (incubating) Impact A stronger, converged, community model

Unity with Community



Dremio is excited to help bring the various functions and capabilities of Nessie into the Apache Polaris (incubating) project.

- Inclusive community
- Robust open source catalog for open lakehouse architectures
- Reduce catalog sprawl
- Broader group of contributors

This partnership not only accelerates technical progress but also brings more contributors into the Nessie community, further strengthening the growing ecosystem around Polaris

Thank You

