



WHITEPAPER

Reducing the Cost of Cloud Data Analytics: 3 Architecture Choices

By Tomer Shiran
Founder and Chief Product Officer, Dremio

The need to do more with less

IT budgets are under pressure due to economic uncertainty, forcing technology leaders to find ways to accomplish more with less. With data and technology at the heart of the business, it is not possible to simply shut down cloud migrations and data analytics projects. Furthermore, in some verticals such as financial and health services, higher volatility and volume is leading to increasing amounts of data that needs to be processed and analyzed.

In this paper we look at three popular architecture choices for cloud data analytics, then describe how Dremio can help you accelerate projects and productivity at a fraction of the cost of cloud data warehouses and simple query engines.

The rise of S3 and ADLS as the new systems of record

Cloud data lake storage (S3 and ADLS) has emerged as the de-facto system of record for data in the cloud. By “system of record” we mean the data store from which all other data stores can be reconstituted. It is the data you cannot be without. There are numerous attributes that led to the rise of S3 and ADLS as the systems of record in the public cloud, including:

- **Infinite scalability.** There are no practical limits to how much data companies can store in S3 and ADLS. Companies can simply save as much data as they want into these services.
- **Low cost.** Cloud data lake storage is inexpensive. The cost is about \$0.20 per TB/month, and you only pay for the actual data stored as opposed to storage provisioned.
- **Availability.** Cloud data lake storage is highly available. For example, according to AWS, the standard tier of S3 is designed for 99.99% availability (less than one hour downtime in a year) and 99.999999999% durability.
- **Near-zero administration.** There's very little administration required with S3 and ADLS compared to traditional storage. You basically interact with a global, highly available service.
- **Global.** S3 and ADLS are available in dozens of regions and countries, making it possible to meet latency and data residency requirements.

For companies migrating to the cloud and adopting cloud-native architectures, data lands in S3 and ADLS in two ways. First, applications that were born in the cloud often default to S3 and ADLS for their data storage. Second, data pipelines are created from on-premise applications and databases in order to move data to S3 and ADLS. Once the required data is in S3 and ADLS, it is then copied and moved into a variety of systems, such as databases, data warehouses, and search engines, in order to meet the needs of different users and applications.

The options for cloud data analytics

Data consumers looking to derive insights and value from the rapidly increasing amount of critical data landing in S3 and ADLS have a few options for enabling data consumers to query it. The following table provides a high-level comparison:

	CLOUD DATA WAREHOUSE	QUERY ENGINE	DATA LAKE ENGINE
Examples	Redshift, Snowflake	Presto, Athena	Dremio
Cost	\$\$\$	\$\$	\$
Data stays in S3 & ADLS	No	Yes	Yes
Interactive BI speed	Yes (extracts may be required)	No	Yes
Public cloud	Yes	Yes	Yes
Hybrid cloud	No	Athena - No Presto - Yes	Yes

CLOUD DATA WAREHOUSE: REDSHIFT & SNOWFLAKE (\$\$\$)

Unlike on-premise data warehouses such as Teradata (also relatively expensive), cloud data warehouses do provide an option to start small at a reasonable cost, but then can become **very** expensive when used in production. Many companies that adopt a cloud data warehouse may later find that the actual costs are far greater than anticipated, and are forced to consider cost control measures. At this point, the common approach is to seek an alternative, more cost-effective model for the majority of their data analytics use cases.

Cloud data warehouses also present a cloud-architecture problem. The way to achieve optimal performance in a data warehouse is by copying and moving data directly into it. This violates the design of modern cloud architectures (aka [Next Architecture](#)) and can introduce significant cost due to:

- The compute cost of ingesting all data into the data warehouse
- The storage cost of saving another copy of all the data in the data warehouse (it's already in your S3/ADLS account)
- The engineering costs associated with building and maintaining a complex data pipeline

QUERY ENGINE: PRESTO & ATHENA (\$\$)

The query engine approach moves closer to cloud architecture principles by ensuring the data stays in one location. This is a step forward from data warehouses, but introduces other significant challenges. Presto is an open source query engine that was created by Facebook, and AWS Athena is a Presto-based service. We will now look at each of these in brief.

Presto is an open-source project and, as such, is considered “free”, but the infrastructure on which it runs is not. Furthermore, because Presto was designed by Facebook to run on limitless clusters of servers, it leaves a lot to be desired in terms of efficiency. As a result, the infrastructure cost of a Presto cluster can quickly become prohibitively high — especially when running in a cloud infrastructure where cost/query has become the most important consideration. More problematic than cost is performance. Presto does not provide the high concurrency and low latency required for self-service BI tools like Tableau, Power BI, Looker and Superset, thereby disqualifying it from the vast majority of use cases for data consumers..

Amazon’s Athena service eliminates the overhead of managing a Presto cluster by routing customer queries to shared Presto clusters. This approach provides a utility pricing model in which you pay per TB scanned. Athena is an inexpensive solution if you run just a few queries per week, but can become impractical for many real-world production workloads. Furthermore, there have been [cases](#) in which mistyped SQL queries or scripts resulted in enormous bills (hundreds of thousands of \$). In addition, because the Presto clusters are shared in the Athena architecture, performance and reliability can vary based on time of day and so-called “noisy neighbors.”

DATA LAKE ENGINE: DREMIO (\$)

Our third architecture option is a data lake engine, which preserves the core principles of cloud architecture while addressing the inherent performance and functionality problems found in simple query engines. A data lake engine combines query acceleration and elasticity to provide significant performance and cost benefits.

Dremio is a data lake engine that enables high-performance queries directly on S3 and ADLS by accelerating query execution. This is accomplished with a number of query acceleration technologies:

- **Apache Arrow-based engine.** Created by Dremio, Apache Arrow is an open source columnar in-memory technology with over 10 million monthly downloads. Dremio’s internal execution engine is based entirely on Arrow, leveraging LLVM-based code generation to achieve 3-4x faster ad-hoc query speed compared to Presto.

Query acceleration
reduces cloud
infrastructure costs by

>75%

- **C3 (Columnar Cloud Cache).** S3 and ADLS are infinitely scalable and cheap, but suffer from much higher latency than local SSD storage. In addition, cloud providers charge ~\$0.40 per million requests, representing 10-20% of the amortized cost of a single query. Dremio includes a patent-pending distributed caching system which leverages the local SSDs on the instances/virtual machines to transparently cache 8KB blocks as they are read from cloud data lake storage. Because most data is accessed hundreds or thousands of times, C3 drastically speeds up I/O and also eliminates S3/ADLS request costs. Note that the Dremio query optimizer is cache-aware, so it makes intelligent decisions on how to distribute the work involved in running a query.
- **Data Reflections.** In many cases it is simply impossible to achieve interactive performance with a full table scan of the raw data. Dremio provides a patent-pending technology called Data Reflections, which accelerates queries by creating and maintaining optimized data structures. The query optimizer automatically rewrites incoming query plans to take advantage of the reflections that are available in the system. Note that the reflections are typically persisted as Apache Parquet files on S3/ADLS and refreshed periodically by the system, although you also have the option to create them independently and simply notify the optimizer with an API call. With reflections in place, BI queries on data lake storage often run 100x faster than Presto, enabling direct access by Tableau, Power BI and other tools without having to create extracts, aggregations or cubes.

These query acceleration technologies combine to deliver 4-100X faster performance compared to Presto. This can also translate directly into significant cost savings, since less infrastructure is required to achieve the same performance. For example, a 4X average speed increase results in a 75% infrastructure cost reduction vs. Presto at the same level of performance.

ELASTIC EXECUTION REDUCES CLOUD INFRASTRUCTURE COSTS BY > 60%

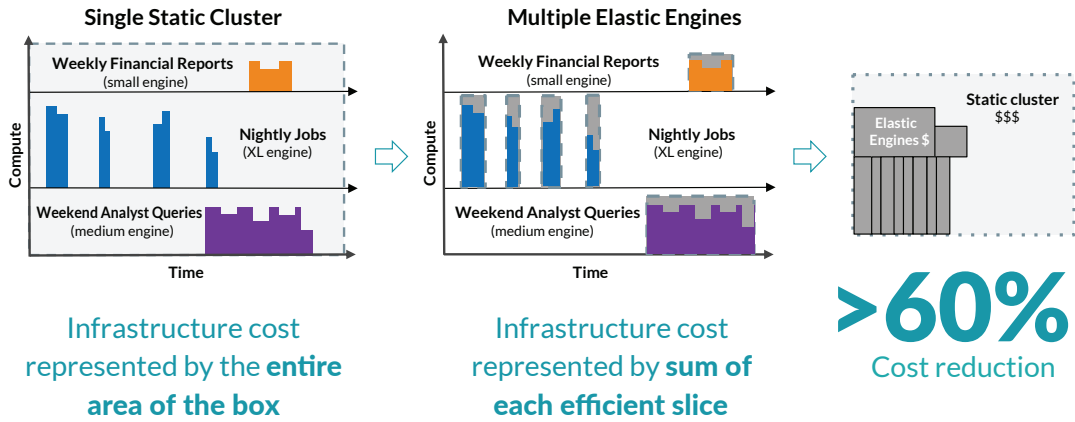
Most distributed systems must be sized based on peak workloads. For example, you would need to operate a 100-node Presto cluster if you expected the workload to require that capacity at any time. This is extremely inefficient, resulting in idle resources most of the time. What if the infrastructure could automatically adjust based on the actual workload? And what if you could also ensure that different workloads were completely isolated from each other? This is now possible with some systems such as Dremio and Snowflake.

Dremio's recently released Elastic Engines feature combines a single control plane with zero to N execution engines. The control plane is responsible for planning queries and routing them to the appropriate engine based on pre-defined rules or manual selection. You can leverage this capability to utilize different engines for different workloads or users, independently sizing each engine for that specific workload. An engine automatically goes to sleep when there are no queries to run.

4-100^x

faster performance
compared to Presto

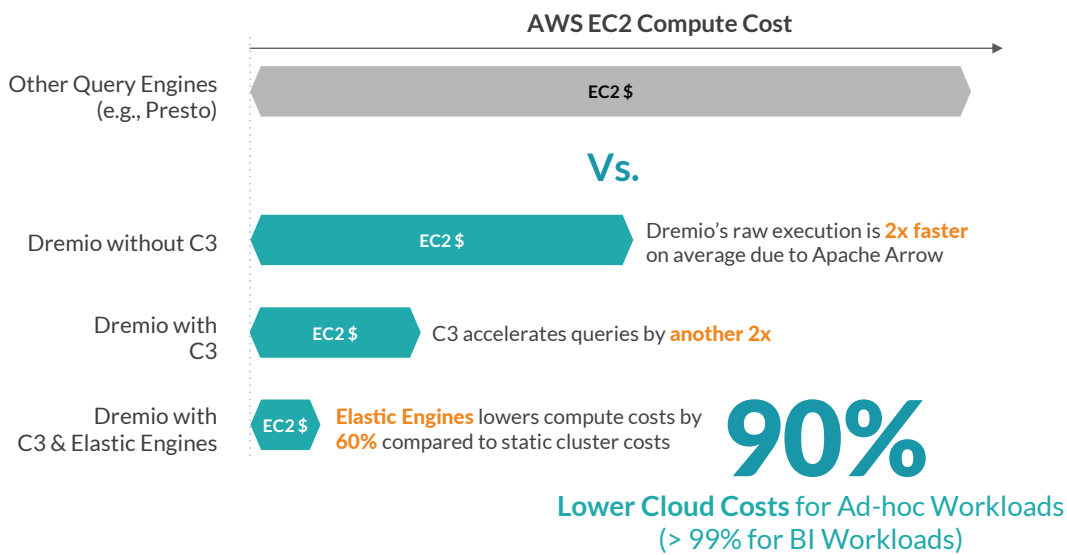
As a result, you can size each engine according to the workload it is serving, as shown in the following illustration:



Based on our experience, the ability to automatically adapt the compute capacity to the workloads leads to a cost reduction of over 60% vs. running on a static cluster of instances/virtual machines.

THE COMPOUNDING COST SAVINGS OF QUERY ACCELERATION & ELASTIC ENGINES

Query Acceleration (Apache Arrow, C3 and Data Reflections) provides a performance increase of 4-100x compared to simple query engines, resulting in a cost reduction of 75% (4x) or more. Elastic Engines provide an additional cost reduction of 60% (2.5x) or more. Combined, Dremio is able to reduce the cloud infrastructure spend for data lake analytics by 90% ($4 * 2.5 = 10$) as shown in this illustration:



What if you still have on-premise data lakes?

The economic slowdown resulting from the Coronavirus pandemic may, unfortunately, also slow down your cloud migration plans. Over the last 10 years, many companies have built on-premise data lakes based on the Apache Hadoop stack (or, in some cases, using S3-compatible data stores like ECS or Minio).

Dremio was designed for the public cloud, but many companies also run it on an on-premise data lake, realizing the exact same advantages — lightning-fast query speed and a self-service semantic layer that empowers business analysts and makes data engineers more efficient. And when you're ready to resume your migration to the cloud, Dremio migrates with you, and you can continue to realize the benefits from your initial time and \$ investments.



ABOUT DREMIO

Dremio delivers lightning-fast queries and a self-service semantic layer directly on your data lake storage. No moving data to proprietary data warehouses, no cubes, no aggregation tables or extracts. Just flexibility and control for data architects, and self-service for data consumers.

Deploy Dremio

[Learn more at dremio.com](https://dremio.com)

CONTACT SALES

contact@dremio.com