# The Power of Presto

Open Source for the Enterprise

Mark Shainman

TERADATA.

When it comes to getting real business value out of big data, we must think about what users want. Do they want more data or less? Do they want queries to be easy or complicated? Do they want to do a lot of work to get an answer or make everything easy? The trick is not just to be able to do a proof of concept with big data, but to make it usable by everyone in your company who needs access to it.

Why is expanding access to big data easier said than done? To date, accessing big data stores has often required specialized skills. Big data teams are sometimes siloed from other analytical teams. Analysts asking questions of big data must interact with these teams to find out about the relevant information in big data and make requests to pull data or perform analyses on their behalf. This complicates access to big data.

Such a process is also inefficient: the questions analysts are posing today often need answers from multiple data stores. For example, information about product sales data comes from traditional sources such as the data warehouse while social media buzz about the product is derived from big data sources. Analysts must manage multiple queries and later join the results together in a meaningful way to find answers to their questions.

## How Can We Make Wider Use of Big Data?

First of all, we must empower analysts to use the skills they have by supporting the lingua franca of analytics: SQL. SQL access to big data breaks down big data silos from a skills perspective.

Second, we must enable interactive, ad-hoc queries. SQL queries on big data have not been impossible in the past, but understanding and using tools such as Hive for handling such queries was difficult for most analysts. Furthermore, even if analysts did use Hive, interactive query performance wasn't the best. Modern analysts don't expect to ask a question today and get the answer tomorrow; they need an interactive tool.

Third, wider use of big data requires seeing big data as one of many data sources. Rather than framing separate queries for separate data stores, analysts should be able to ask questions that reach out to multiple sources and platforms and bring back answers as needed to answer the query. This begins breaking down the siloing of big data from other data repositories.

An even greater expansion of the use of big data requires enabling users who may or may not know SQL themselves to leverage familiar BI tools to generate these queries behind the scenes. This entails not just supporting SQL, but the machine generated SQL that BI tools speak.

TERADATA

It requires robust SQL support that doesn't rely on human-crafted queries. BI tools must form the query from their generated SQL and return a response that is digestible by these tools.

With all of these capabilities in place, we can envision broader, more integrated access to data in general and big data in particular. We can see it leveraged by many more people across the organization, not just those with specialized skills in the many tools that comprise the Hadoop ecosystem.

Thankfully, a great deal of the work to meet the requirements just described has already been done and made available as open source.

## Enter Facebook's Presto

Facebook recognized the need for a higher performing, interactive SQL on Hadoop solution and, in response, developed Presto. Presto is a powerful, next-generation, open source SQL query engine that supports big data analytics. Designed for low-latency interactive data analysis, Presto is best suited for workloads that require a faster query engine to support interactive speeds for data exploration as well as a wide variety of connectors to query multiple data sources.

Presto is designed to be faster than Facebook's earlier Hadoop data query framework, Hive, for interactive workloads on large-scale data. And it is. Presto's full memory-based architecture enables it to run magnitudes

faster than traditional Hive. As a result, Presto has become a key component of Facebook's analytical infrastructure – the company runs tens of thousands of queries a day on data stores that scale up to 300 petabytes. Facebook released Presto as open source in 2013 and today it is used by many large organizations, including Netflix, Dropbox, Airbnb, and NASDAQ.

## The Need for a Full Toolset

Does Presto replace Hive? No. Hive is widely used and is particularly strong for its ability to restart large scale jobs if they are interrupted as well as to perform large-scale data analysis with many joins. For a detailed comparison of Presto and Hive, see "Presto: A More Scalable Flexible Open Source SQL on Hadoop Engine."

Although this paper describes Presto, it's important to bear in mind that Presto is one tool among many. The vision of access to all the data by all authorized users has many moving parts and realizing that vision requires different tools that accomplish different business purposes. Organizations may leverage strong tools associated with their Hadoop distribution as well as selecting tools from among the many active open source projects.
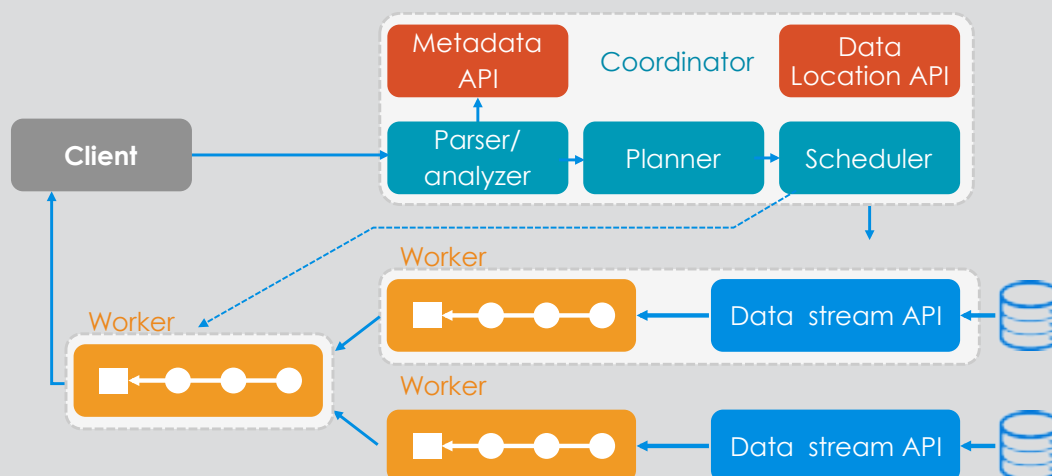


Figure 1: Presto Architecture

TERADATA.

## Presto Key Features

**Distribution-agnostic.** Presto runs on multiple Hadoop distributions, making it one of the most portable SQL-on-Hadoop tools. This is particularly important because everyone does not use the same Hadoop distribution.

**Multiple data sources**. Presto is not limited to querying data stored in Hadoop; it allows users to take analytics, power, and data processing to wherever the data resides. Presto has a connector based architecture that enables it to reach out to numerous types of data stores and platforms. Leveraging this connecter architecture, Presto can, through a single query, access data that resides in Hadoop, NoSQL databases, relational databases, and proprietary data stores. For example, Presto can reach out to data in PostgreSQL, MySQL and Cassandra.

**No MapReduce**. Presto doesn't use MapReduce. It's a SQL engine that directly accesses data in Hadoop. As a result, it's much faster than traditional Hive for large, interactive workloads.[1] On the other hand, Hive is stronger for long-running, complex jobs.

**Scalability.** Scalability is a critical concern, because effective data use increases the number of people using it, both directly and through BI tools. Presto delivers scalability

due to its modern code base. What's more, this scalability is proven in enterprise production environments, like Facebook. On average, Netflix runs about 3,500 Presto queries on its Hadoop clusters every day.

## Benefits of Presto and Teradata

In *Open Source for the Enterprise* (O'Reilly), Dan Woods speaks about needing to productize open source software. Facebook recognized the growth in the Presto community in a blog post, citing use of and work on Presto by large organizations including (among others) Airbnb, Dropbox, Netflix, NASDAQ, Microstrategy, and Teradata, as well as international companies such as Japanese social media game-development company Gree and JD.com, a Chinese ecommerce company.

Teradata is strongly supporting Presto development. In fact, Teradata has made a multi-year commitment to contributing open source development to Presto, with more than 20 full time employees dedicated to contributing code for the project. The vast majority of the work is 100% open source and is free to download. This initiative is in large part the result of Teradata's acquisition of Hadapt, a pioneer in SQL on Hadoop.
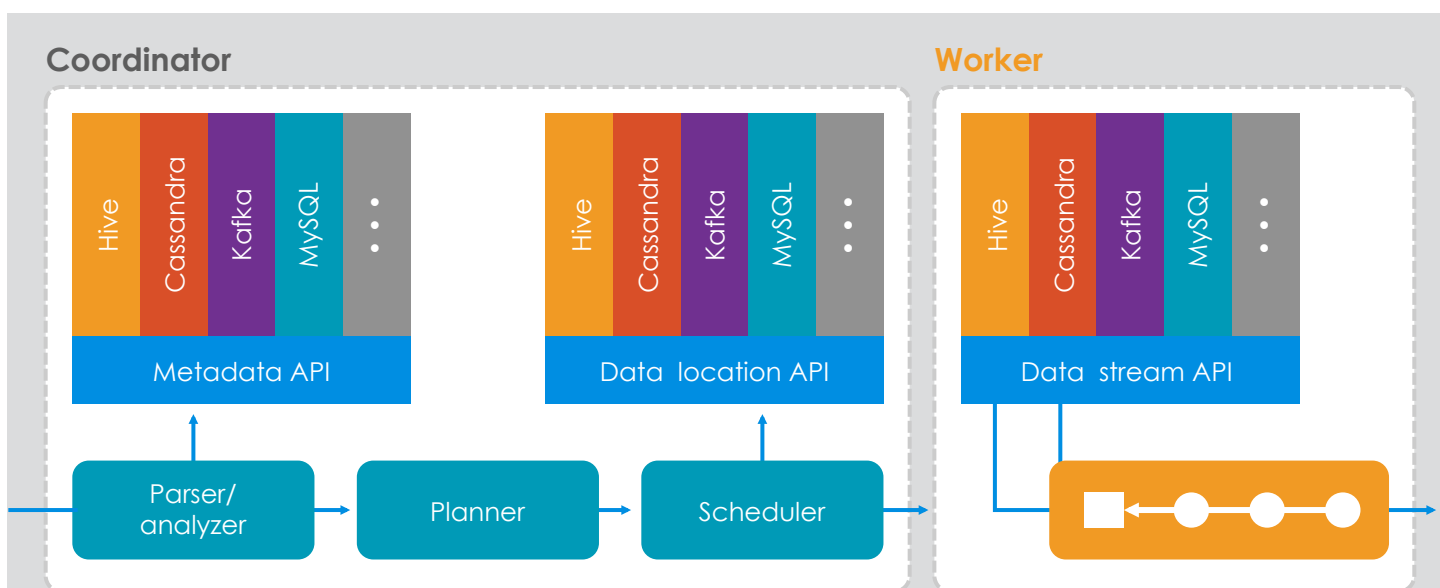


Figure 2: Presto extensibility - connector interfaces

**TERADATA**

Teradata's involvement with Presto has centered around the fit and finish of Presto. The work has focused on the advancement of the functionality, ease of use and security of the core Presto SQL engine. The company has provided certified, stable releases of open source Presto as well as add-ons and connectors. The goal is to enable Presto to work simply and seamlessly within the enterprise.

Facebook recently recognized Teradata's contribution, saying, "A special shout out goes to Teradata — which joined the Presto community this year with a focus on enhancing enterprise features and providing support — for having seven of our top 10 external contributors."

– Open source in 2015:
   A year of growth

## Examples of Teradata's Contributions

- Full documentation for Presto

- Presto Admin, a tool for installing, managing, and configuring Presto on a Hadoop cluster

- Enterprise class ODBC and JDBC drivers, developed on an accelerated schedule in response to strong demand from the community.

- A Presto plug-in for Ambari, an open source tool for managing Hadoop via a Web UI

- YARN integration, to enable resource management for Presto

- RPMs to make it easier to install Presto

Teradata also offers enterprise support for Presto. The technical support comes from a company that is steeped in Presto and has a deep understanding of the code. Teradata has also been working with Facebook on a roadmap for its developments, which provides the community with a strong sense of what is coming and eases enterprise adoption.

## Presto's Role in the Unified Data Architecture

Most companies are deploying an ecosystem of analytical engines, and Hadoop is an important and growing part of that ecosystem. This ecosystem approach was a driving force behind the creation of Teradata's Unified Data Architecture (UDA). Teradata embraces an ecosystem because no single tool can service all needs. Facebook concluded this many years ago, and it drove the engineering around Presto. But ecosystems also have challenges such as different protocols for different engines, and being able to author queries that reach into multiple systems. Presto addresses these challenges.

So where does Presto fit within the UDA? Simply put, Presto enables the UDA to work better. Integration across the entire UDA is crucial, and Presto is a natural fit to do just that. Presto links the data warehouse and data lake, and allows the data lake to participate in integrated analytics, thus expanding the scope of the data platforms that the UDA can touch. By purchasing and using the Teradata QueryGrid connector for Presto, users can execute Presto queries from Teradata to Hadoop with pushdown processing into Hadoop. Users can also execute Presto queries on Hadoop that query data in Teradata with pushdown processing into Teradata.

Presto enables a flexible architecture of effective analytics across multiple platforms. It provides a component of the glue to keep all that data together, while also improving the engines that allow users to access the data that lives in Hadoop and other data stores. With Presto, the UDA works better and is more extensive. Now with Presto, the UDA not only has the ability to query Hadoop data in an interactive manner, but also to access data from other data sources that were not formally directly queryable from within the UDA, including platforms such as MySQL, PostgreSQL, Cassandra, and Kafka.

Presto's ability to leverage ANSI SQL is key because it enables SQL to become the common language across all the access points within the UDA. Queries can come from SQL literate analysts and, more broadly, any BI tool that speaks SQL (which the vast majority of BI tools do).

TERADATA.

Presto can also be used to take a huge amount of data, query it, create an aggregate set and then join it within a Teradata environment. This is cost efficient and production effective.

## Conclusion

Presto addresses a real need for a portable SQL on Hadoop tool. It is architected from the ground up for high performance interactive query processing. Open source is a fount of continual innovation, especially with regard to big data. In addition, there are strong tools that come with specific Hadoop distributions.

The fact is that organizations will deploy multiple tools. Presto has many strengths, including full SQL support, interactive queries, the ability to reach out to multiple data sources with a single query, proven scalability, and low barrier to entry in terms of costs and skills. Presto is therefore worthy of consideration as part of the enterprise toolset that will empower your organization to make effective use of all its data, including data stored in Hadoop. For organizations moving toward a Unified Data Architecture, the rationale for adopting Presto is even stronger.

To learn more or try out Presto, visit teradata.com/ PrestoDownload.

## End Notes

1. For information on how Netflix uses both Presto and Hive, see http://techblog.netflix.com/2014/10/using-presto-in-our-big-data-platform.html.

TERADATA.