

<https://www.qubole.com/>[Home...](#) (<https://www.qubole.com/>) > [Blog](#) (<https://www.qubole.com/blog>) > Top Apache Spark Use Cases

Blog

Top Apache Spark Use Cases

By Ari Amster

March 10, 2016



(<https://4cti9u1qldo011vlp4eutwnxf4-wpengine.netdna-ssl.com/wp-content/uploads/2016/03/810x430-top-apache-cases-expanded.jpg>)

This post was originally published in July 2015 and has since been expanded and updated.

Apache Spark is quickly gaining steam both in the headlines and real-world adoption. UC Berkeley's AMPLab developed Spark in 2009 and open sourced it in 2010. Since then, it has grown to become one of the largest open source communities in big data with over 200 contributors from more than 50 organizations. This open source analytics engine stands out for its ability to process large volumes of data significantly faster than MapReduce because data is persisted in-memory on Spark's own processing framework.

When considering the various engines within the Hadoop ecosystem, it's important to understand that each engine works best for certain use cases, and a business will likely need to use a combination of tools to meet every desired use case. That being said, here's a review of some of the top use cases for Apache Spark (<http://www.qubole.com/apache-spark-as-a-service/>).

1. Streaming Data

Apache Spark's key use case is its ability to process streaming data. With so much data being processed on a daily basis, it has become essential for companies to be able to stream and analyze it all in real time (<http://www.qubole.com/blog/big-data/real-time-analytics/>). And Spark Streaming has the capability to handle this extra workload. Some experts even theorize that Spark could become the go-to platform for stream-computing applications, no matter the type. The reason for this claim is that Spark Streaming unifies disparate data processing capabilities, allowing developers to use a single framework to accommodate all their processing needs.

Among the general ways that Spark Streaming is being used by businesses today are:

Streaming ETL – Traditional ETL (extract, transform, load) tools used for batch processing in data warehouse environments must read data, convert it to a database compatible format, and then write it to the target database. With Streaming ETL, data is continually cleaned and aggregated before it is pushed into data stores.

Data enrichment – This Spark Streaming capability enriches live data by combining it with static data, thus allowing organizations to conduct more complete real-time data analysis. Online advertisers use data enrichment to combine historical customer data with live customer behavior data and deliver more personalized and targeted ads in real-time and in context with what customers are doing.

Trigger event detection – Spark Streaming allows organizations to detect and respond quickly to rare or unusual behaviors (“trigger events”) that could indicate a potentially serious problem within the system. Financial institutions use triggers to detect fraudulent transactions and stop fraud in its tracks. Hospitals also use triggers to detect potentially dangerous health changes while monitoring patient vital signs—sending automatic alerts to the right caregivers who can then take immediate and appropriate action.

Complex session analysis – Using Spark Streaming, events relating to live sessions—such as user activity after logging into a website or application—can be grouped together and quickly analyzed. Session information can also be used to continuously update machine learning models. Companies such as Netflix use this functionality to gain immediate insights as to how users are engaging on their site and provide more real-time movie recommendations.

2. Machine Learning

Another of the many Apache Spark use cases is its machine learning capabilities.

Spark comes with an integrated framework for performing advanced analytics that helps users run repeated queries on sets of data—which essentially amounts to processing machine learning algorithms. Among the components found in this framework is Spark’s scalable Machine Learning Library (MLlib). The MLlib can work in areas such as clustering, classification, and dimensionality reduction, among many others. All this enables Spark to be used for some very common big data functions, like predictive intelligence, customer segmentation for marketing purposes (<https://www.qubole.com/resources/webinars/big-data-analytics-marketers/>), and sentiment analysis (<https://www.qubole.com/resources/webinars/sentiment-analysis-using-hive-secrets-pros/>). Companies that use a recommendation engine will find that Spark gets the job done fast.

Network security is a good business case for Spark’s machine learning capabilities (<https://www.qubole.com/resources/webinars/machine-learning-deployment/>). Utilizing various components of the Spark stack, security providers can conduct real time inspections of data packets for traces of malicious activity. At the front end, Spark Streaming allows security analysts to check against known threats prior to passing the packets on to the storage platform. Upon arrival in storage, the packets undergo further analysis via other stack components such as MLlib. Thus security providers can learn about new threats as they evolve—staying ahead of hackers while protecting their clients in real time.

3. Interactive Analysis

Among Spark’s most notable features is its capability for interactive analytics. MapReduce was built to handle batch processing, and SQL-on-Hadoop engines such as Hive or Pig are frequently too slow for interactive analysis. However, Apache Spark, is fast enough to perform exploratory queries without sampling. Spark also interfaces with a number of development languages including SQL, R, and Python. By combining Spark with visualization tools, complex data sets can be processed and visualized interactively.

Debuting in April or May of this year, the next version of Apache Spark (Spark 2.0) will have a new feature—*Structured Streaming*—that will give users the ability to perform interactive queries against live data. Combining live streaming with other types of data analysis, Structured Streaming is predicted to provide a boost to Web analytics (<https://www.qubole.com/resources/webinars/using-big-data-tools-analyze-web-analytics-data/>) by allowing users to run interactive queries against a Web visitors current session. It could also be used to apply machine learning algorithms to live data. In this scenario the algorithms would be trained on old data and then redirected to incorporate new—and potentially learn from it—as it enters the memory.

4. Fog Computing

While big data analytics may be getting a lot of attention, the concept that really sparks the tech community’s imagination is the Internet of Things (http://go.qubole.com/2016-06-14---WR---IoT-Case-Study-Using-Big-Data-in-the-Cloud-wGCP_LP.html) (IoT). The IoT embeds objects and devices with tiny sensors that communicate with each other and the user, creating a fully interconnected world. This world collects massive amounts of data, processes it, and delivers revolutionary new features and applications for people to use in their everyday lives. However, as the IoT expands so too does the need for distributed massively parallel processing of vast amounts and varieties of machine and sensor data. All that processing, however, is tough to manage with the current analytics capabilities in the cloud.

That’s where fog computing (<http://www.informationweek.com/big-data/big-data-analytics/apache-spark-3-promising-use-cases/a/d-id/1319660>) and Apache Spark come in.

Fog computing decentralizes data processing and storage, instead performing those functions on the edge of the network. However, Fog computing brings new complexities to processing decentralized data, because it increasingly requires low latency, massively parallel processing of machine learning, and extremely complex graph analytics algorithms. Fortunately, with key stack components such as Spark Streaming, an interactive real-time query tool (Shark), a machine learning library (MLlib), and a graph analysis engine (GraphX), Spark more than qualifies as a fog computing solution. In fact, as the IoT industry gradually and inevitably converges, many industry experts predict that—compared to other open source platforms— Spark has the potential to emerge as the de facto fog infrastructure.

Spark in the Real World

As mentioned earlier, online advertisers and companies such as Netflix are leveraging Spark for insights and competitive advantage. Other notable businesses also benefitting from Spark are:

Uber – Every day this multinational online taxi dispatch company gathers terabytes of event data from its mobile users. By using Kafka, Spark Streaming, and HDFS, to build a continuous ETL pipeline, Uber can convert raw unstructured event data into structured data as it is collected, and then use it for further and more complex analytics.

Pinterest – Through a similar ETL pipeline, Pinterest can leverage Spark Streaming to gain immediate insight into how users all over the world are engaging with Pins—in real time. As a result, Pinterest can make more relevant recommendations as people navigate the site and see related Pins to help them select recipes, determine which products to buy, or plan trips to various destinations.

Conviva – Averaging about 4 million video feeds per month, this streaming video company is second only to YouTube. Conviva uses Spark to reduce customer churn by optimizing video streams and managing live video traffic—thus maintaining a consistently smooth, high quality viewing experience.

When NOT to Use Spark

Even though it is versatile, that doesn't necessarily mean Apache Spark's in-memory capabilities are the best fit for all use cases. More specifically, Spark was not designed as a multi-user environment (<http://www.cio.com/article/2935621/big-data/ibm-commits-to-apache-spark-compute-engine.html>). Spark users are required to know whether the memory they have access to is sufficient for a dataset. Adding more users further complicates this since the users will have to coordinate memory usage to run projects concurrently. Due to this inability to handle this type of concurrency, users will want to consider an alternate engine, such as Apache Hive, for large, batch projects.

Over time, Apache Spark will continue to develop its own ecosystem, becoming even more versatile than before. In a world where big data has become the norm, organizations will need to find the best way to utilize it. As seen from these Apache Spark use cases, there will be many opportunities in the coming years to see how powerful Spark truly is.

As more and more organizations recognize the benefits of moving from batch processing to real time data analysis, Apache Spark is positioned to experience wide and rapid adoption across a vast array of industries

Interested in learning more about Apache Spark, collaboration tools offered with QDS for Spark, or giving it a test drive? Click the button to learn more about Apache Spark-as-a-Service.

LEARN ABOUT SPARK ([HTTP://WWW.QUBOLE.COM/APACHE-SPARK-AS-A-SER](http://www.qubole.com/apache-spark-as-a-ser))

ONE THOUGHT ON "TOP APACHE SPARK USE CASES"



Kumar

January 22, 2017 at 12:23 pm (<https://www.qubole.com/blog/big-data/apache-spark-use-cases/#comment-7618>)

informative and a quick overview of what spark is capable of.

REPLY ([HTTPS://WWW.QUBOLE.COM/BLOG/BIG-DATA/APACHE-SPARK-USE-CASES?REPLYTOCOM=7618#RESPOND](https://www.qubole.com/blog/big-data/apache-spark-use-cases?replytocom=7618#RESPOND))

LEAVE A REPLY

Your email address will not be published. Required fields are marked *

Comment

Name *

Email *

Website

Are you human? *

seven + = sixteen ↻

Post Comment

Live Session



(<http://bit.ly/2cbRkXh>)

On-demand Webinar
