

Broad Institute's Use of Google Cloud: Accelerating Biomedical Research for All

Google engineering, Intel optimizations, and Intel® Xeon® Scalable processor-based cloud instances drive down costs and enable new levels of research and collaboration

Broad Institute's Google Cloud solution:

- Google Cloud
- The Broad Institute Genome Analysis Toolkit (GATK)
- Intel® Xeon® Scalable processors
- Intel® Genomics Kernel Library
- Intel® Advanced Vector Extensions 512
- Intel® Intelligent Storage Acceleration Library



Executive Summary

[The Broad Institute of MIT and Harvard \(Broad Institute\)](#) is a world-renowned nonprofit academic organization dedicated to transforming medicine by

- Discovering the molecular basis of major human diseases.
- Developing effective new approaches to diagnostics and therapeutics.
- Openly disseminating discoveries, tools, methods, and data to the entire scientific community.

To support these objectives, the Broad Institute uses an array of cutting-edge technologies, including genome sequencing and analysis, which produce very large datasets.

The Broad Institute began migrating its genomic data storage and analysis workloads to Google Cloud in 2014. Increases of data volume had been fueled by falling costs of data generation, as well as a wave of scientific and technological innovation. Using the cloud would set a foundation for a sustainable future of genomics research. The cloud would deliver scalable computing resources as well as data management facilities and tools the Broad Institute needed for growth. The cloud vision involved workload migration and building a new collaboration platform for data sharing, analysis and collaboration called Terra.

In collaboration with Google Cloud and Intel, the Broad Institute optimized their genomics workloads for fast, cost-effective execution on Google Cloud N1 and N2 instances. Compared to the initial deployment of workloads on Google Cloud, the collaboration resulted in 85 percent reduction in cost of data processing after optimization.¹

Terra is co-developed through a partnership with Microsoft and Verily. Terra packages the relevant functionality and other cloud-based resources in a form that is more immediately usable and broadly accessible by researchers in the life sciences. It is built on top of the cloud infrastructure, aiming to enable next-generation collaborative research through an open, scalable, and secure platform. Terra allows researchers to access data, run analysis tools, and work together, giving them resources and an environment with which they can achieve scientific breakthroughs. The platform enables the Broad Institute to make its optimized tools and workflows openly accessible to researchers around the world, maximizing the impact of its research and investments in technology development.

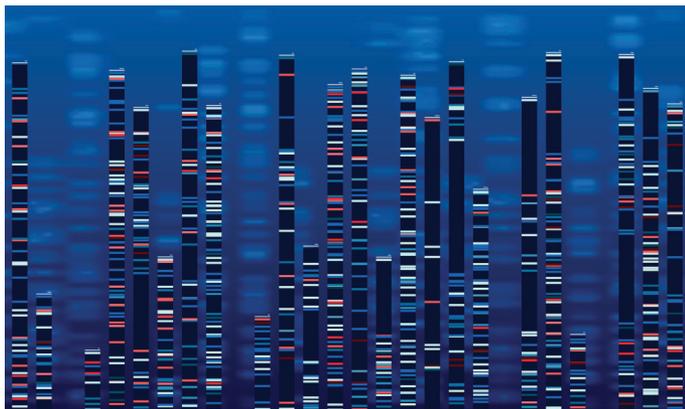
Challenge

Modern life sciences research uses genome sequencing and genomic analysis. These help scientists access the roots of human biology, understand the many challenges to human health, and create new treatments and therapies to battle diseases. Analyzing genomic data involves sophisticated computational tools, such as the Genome Analysis Toolkit (GATK), which the Broad Institute develops in-house. The organization makes the GATK available to researchers around the world in the form of open-source packages, as well as fully-assembled genomic analysis pipelines called GATK Best Practices workflows.

The Broad Institute is a major center of genomics-based research, not only as a provider of analytics tools and computational resources, but also as a sequencing service provider.

“Our genome sequencing facility generates data on the order of a whole human genome about every three to five minutes, 24 hours a day,” Geraldine Van der Auwera, Director of Outreach and Communications at the Broad Institute's Data Sciences Platform, explained. “Each genome corresponds to approximately 350 gigabytes of data before compression, resulting in some 30 petabytes so far of genomic data being managed by the Broad Institute.”

Like many cutting-edge technologies, genome sequencing started out being very expensive—the seminal Human Genome Project cost USD 2.7 billion. Over the last decade, the cost of sequencing has decreased dramatically to less than a thousand dollars per whole genome today. Lower cost makes this valuable tool more accessible to researchers and has resulted in a flood of genomic data.



To support their objective to transform medicine, the Broad Institute uses an array of cutting edge technologies, including genome sequencing and analysis

Access to more data has opened doors for innovative and expanded research. But, larger and larger datasets require ever-expanding storage capacity and computational capability. At the Broad Institute, data growth placed an increasing amount of strain on its on-premises IT facilities. Compounding this, demand for the institute's computing resources tended to be unevenly distributed throughout the year. High-intensity peaks at certain times put more pressure on their IT systems and staff to meet all researchers' needs in a timely manner. Broad was facing a scalability challenge.

“We realized that the on-premises infrastructure was going to quickly run out of capacity for both storage and computing,” Van der Auwera added. “We decided to go to the cloud for several reasons. One was logistics and economics of operating our processing pipelines and data storage. We could scale as needed for both compute and storage, paying only for the capacity we used, instead of provisioning a data center for seasonal peak workloads. Also, the cloud would allow a whole new level of data federation and collaboration. We could work with others to create a cloud-based data ecosystem, where researchers could combine the data they generated with other datasets into richer, more powerful computational experiments. This would help them achieve greater statistical confidence, integrate additional sources of information, and generate critical insight into the areas of research they were focused on.”

While opening new opportunities, migrating to the cloud presented new challenges.

“We couldn't just copy our existing pipelines over to the cloud,” Van der Auwera commented. “The infrastructures are different. We needed to re-implement our pipelines in a cloud-native way. Plus, to realize our vision of a federated data ecosystem would require building a whole new platform to handle the complexities of the cloud infrastructure, and provide applications and interfaces tailored to the needs of life scientists, in order to enable them to work effectively in the cloud.”

The Broad Institute recognized early on that they would need the expertise of those offering the cloud services for such an undertaking to be successful. While exploring deployment and development options, Google offered to contribute to the development process and helped build the infrastructure needed to deliver on the vision.

“This early collaboration was a key piece of how we began the migration of the institute's production pipelines and setting the foundation for building what would ultimately become the Terra platform,” Van der Auwera said.

Solution

Cloud-native and traditional on-premises cluster computing present vastly different environments. The Broad Institute would have to adapt their pipelines and build a new platform from the ground up that could take advantage of what cloud had to offer.

Leveraging cloud capabilities

“With on-premises infrastructure, you don't have access to a variety of machine types like you can get on the cloud,” Van der Auwera explained. “On-premises clusters are typically all one type of system. With different types of cloud instances, however, we could modularize our workflows and right-size the instances allocated for each task based on its needs. Thus, we could cut processing costs considerably.”

Right-sizing instances is an important aspect of optimizing workflows to run efficiently on cloud.

“Many customers who deploy genomics workflows on cloud reserve large instances, because some parts of the workflows are compute-intensive,” Marissa Powers, an Intel Solution Architect who works with the Broad Institute's data engineering team explained. “The Broad Institute does have

processes in their pipelines that need massive amounts of computation. But most of the tools that are part of the genome analysis pipeline are actually single-threaded. They just need to run as long as they take, and they could use a smaller, less costly instance. So, the Broad Institute team built a sophisticated workflow automation mechanism where individual VMs are right-sized for the job and orchestrated across the entire pipeline of tasks.”

Another key innovation was how they moved data—or rather, how they avoided moving it, whenever possible. Most analysis tools normally require localizing the entirety of the input files from object storage to a VM. The Broad Institute’s GATK can stream just a subset of the genomic data from the original input files. For many stages in the pipeline, execution is parallelized over subsets of the genome, with each subset being sent for processing to a different VM. This streaming approach reduces the amount of storage and memory needed, reduces time spent copying large amounts of data from object storage to the VM, and ultimately reduces costs. Compared to their initial deployment on cloud, these optimizations, along with the use of preemptible instances, reduced the cost of their main genome analysis pipeline by about 85 percent.¹

Optimizing workloads for Intel® architecture

Intel has had a joint partnership with the Broad Institute since 2017, helping optimize the Broad Institute’s pipelines and GATK with Intel libraries, including the Intel Genomics Kernel Library. Intel and the Broad Institute have also collaborated on powerful and flexible data center solutions for genomics analytics for several years. Together they manage the Intel-Broad Center for Genomic Data Engineering. The project optimizes best practices in hardware and software for genome analytics. The Center helps researchers and software engineers build, optimize, and widely share new tools and infrastructure that will help scientists integrate and process genomic data.

Intel worked with the Broad Institute to help optimize their pipelines on Google Cloud. For example, specific kernels in the GATK are optimized for vector operations with Intel Advanced Vector Extensions 512 (Intel AVX-512). Some optimized storage functions use the Intel Intelligent Storage Acceleration Library (Intel ISA-L).

“One kernel of the pipeline, called PairHMM, is a hidden Markov model,” Powers explained. “Intel AVX-512 is a good fit for it based on the length of the vectors being processed. With the optimized version, we’ve seen continuous improvement from the original Java implementation to Intel AVX2 and Intel AVX-512. Anyone who runs the GATK pipeline on Intel Xeon® Scalable processors, gets by default the optimized version, whether they run on-premises or in the cloud.”

Intel benchmarked the Broad Institute’s workloads targeted for Google Cloud. Intel prescribed descriptions for cloud using the Workflow Description Language (WDL), an open-source, community-based standard for pipeline development stewarded by the OpenWDL organization.

The Broad Institute’s pipelines are freely available to anyone to download through [GitHub](#), and to run either on their on-premises infrastructure or the cloud infrastructure of their choice.

Some institutions need to continue to run their workloads and the GATK pipelines on-premises. Intel provides reference architectures for deploying the Broad Institute’s GATK Best Practices workflows to on-premises High Performance Computing infrastructure built on Intel technologies.

Faster and lower cost processing on Google Cloud N2 instances

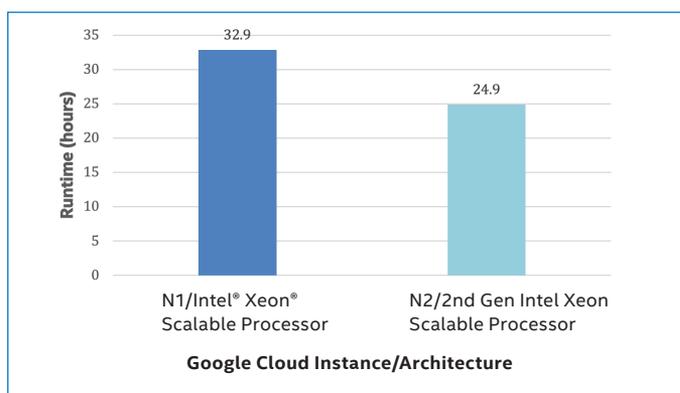
Google Cloud’s N2 instances are built on 2nd Gen Intel Xeon Scalable Processors. Intel benchmarking illustrates how the optimizations of the Broad Institute’s pipelines using Intel AVX-512 and the Intel ISA-L on N2 instances can further accelerate the workloads compared to running on N1 instances (see the following charts).²

As a result, users can run their genomics workflows on Google Cloud about 25 percent faster and at 34 percent lower cost by deploying on N2 instances with 2nd Generation Xeon Scalable processors.²

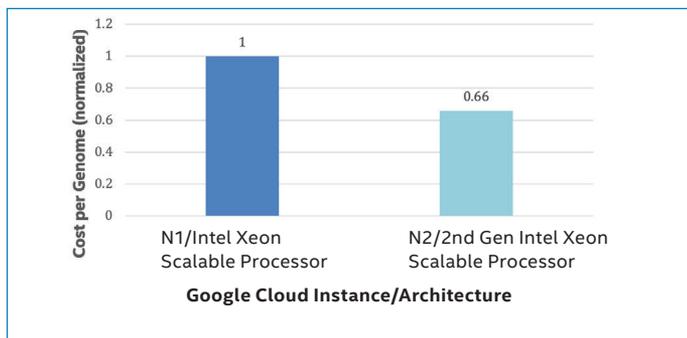
Enabling a Life Sciences research ecosystem with Terra

Genome sequencing, high-resolution medical imaging, and the digital transformation of clinical data have created a sea change in biomedical research. Data science and bioinformatics will reveal new insights out of this data with the help of the right tools, pipelines, and emerging and ongoing research.

The Broad Institute, in collaboration with other organizations and academic institutions, developed a vision of a federated



Pipeline instance runtimes on Google Cloud comparing the first and second generations of Intel Xeon Scalable processors (lower is better)



Instance cost per genome on Google Cloud comparing the first and second generations of Intel Xeon Scalable processors (lower is better)

data ecosystem. This ecosystem would leverage connections between interoperable data repositories, tool repositories, and analysis engines, with user portals tailored to the needs of specific research communities.

Building on the work done to migrate their pipelines to Google Cloud, the Broad Institute partnered with Microsoft and Verily to co-develop Terra. The vision of Terra would advance biomedical research, and it puts powerful tools in the hands of the wider life sciences research community.

"Terra provides a user-friendly environment that enables researchers to access the datasets they need, and apply the tools they want, securely and at scale," Van der Auwera commented. "The platform also makes it easy to share their work at any stage, either privately with their collaborators or publicly with the world, in a form that makes their analysis completely reproducible and extensible."

Terra's built-in emphasis on collaboration allows researchers to tackle human health challenges that are larger than a single organization can solve. Able to draw on a large variety of data, scientists can use established bioinformatics and emerging artificial intelligence (AI) techniques and tools to gain new insights. And thanks to its user-driven design, Terra allows researchers to focus on their science instead of the underlying infrastructure.

Result

By migrating to the cloud and optimizing their workloads, the Broad Institute solved their storage capacity and computational capability challenges in a scalable, forward-looking way. Co-building the Terra platform in partnership with Microsoft and Verily further enabled the Broad Institute to empower not only their own research teams, but life scientists around the world. Terra allows researchers to take advantage of these optimized tools and pipelines and to participate in a federated data ecosystem that opens many exciting new possibilities for biomedical research.



Sequencing and analysis of a genome, shown here as imagined by an artist, generated huge datasets.

According to Van der Auwera, researchers who came from an on-premises infrastructure, where their jobs might often wait in long queues, appreciate the scalability and quick job deployment of Terra. They also enjoy greater flexibility. They can use preloaded workflows, such as the Broad Institute's optimized genomics pipelines, upload their own, or import publicly shared workflows from connected repositories. Dockstore, which is operated by the University of California, Santa Cruz (UCSC), provides such a repository. Bioinformatics scientists who maintain their tools and workflows in Github can register them in Dockstore for use by other researchers, who can then run them on a range of connected analysis platforms including Terra.

"This shows the power of having an ecosystem of platforms that talk to each other," Van der Auwera stated. "It gives researchers enormous flexibility to choose a computational platform that suits their needs and preferences and allows them to collaborate with minimal friction with others across platforms."

Van der Auwera added that many researchers also appreciate Terra's cloud environments framework. The framework makes it possible to run interactive analysis applications, such as Jupyter Notebook, Rstudio, and Galaxy. These run in private, self-contained cloud environments that are preconfigured and can be launched with just a few clicks.

"Currently, we have about 20,000 registered users on Terra on Google Cloud," Van der Auwera concluded, "and it expands steadily with support for Microsoft's Azure cloud currently in the works. And with the data federation, secure sharing, and collaborative capabilities of the cloud becoming more obvious every day, the platform is attracting a multitude of people and organizations with different roles to play in the ecosystem besides researchers themselves. For example, computational tool developers can use it as a convenient method to make their tools available in a way that works out of the box and reduces their support burden. Funding

Case Study | Broad Institute's Use of Google Cloud: Accelerating Biomedical Research for All

agencies and data generators are also very interested in this model as a way to make large data resources widely available without having to maintain their own infrastructure.”

Terra already supports multiple scientific consortia and federal infrastructure development projects. Some of these include:

- The NHGRI's Analysis, Visualization and Informatics Lab-space (AnVIL) project supporting the genomics research community.
- The NCI Cloud Resources supporting the cancer research community.
- The NIH's All of Us Research Program, which aims to gather and process genomic, healthcare, and real-time lifestyle data from 1,000,000 Americans to “learn how our biology, lifestyle, and environment affect health.”³

Solution Summary

To adapt to a dramatic increase in genomics data generation and computational research demand, the Broad Institute migrated their workloads to Google Cloud N1 instances. By modularizing their pipeline workflows, right-sizing cloud instances based on the needs of the workload, and optimizing for Intel Xeon Scalable processors, they reduced processing costs per genome significantly. Running optimized Broad Institute pipelines on Google Cloud N2 instances further accelerates the workloads and reduces costs.

Seeking to realize a wider Life Sciences ecosystem vision, the Broad Institute, Microsoft, and Verily co-developed the Terra platform. Terra is an open, scalable, and secure platform for biomedical researchers to access data, run analysis tools, and collaborate. The platform aims to enable the next generation of collaborative biomedical research by connecting researchers to each other and to the datasets and tools they need to achieve scientific breakthroughs. The Broad Institute makes their optimized pipelines available on Terra, so any researcher in the world can easily run scalable analysis pipelines that are fast, cost-effective, and scientifically excellent.

Running on Google Cloud allows the Broad Institute to scale easily and empower the research community with new capabilities for the benefit of research into solutions to human disease.

Where to Get More Information

Learn more about the [Broad Institute](#) and [Terra](#).

Find out more about [Google Cloud](#)

Explore the capabilities of the [2nd Generation Intel Xeon Scalable processors](#) with integrated Intel Deep Learning Boost capabilities for accelerated AI inferencing.



¹Results reported by the Broad Institute

²Configurations. N1: Testing completed by Intel on 9/10/2020. Region: Google Cloud us-central1 region. N1-standard-2 (2 vCPUs); N1-standard-16 (16 vCPUs). 180+ instances. Intel Xeon Scalable processor (Skylake); memory/instance: 7.5 GB (N1-standard-2); 60 GB (N1-standard-16). Network-attached storage; network bandwidth/instance: 10 Gbps (N1-standard-2); 32 Gbps (N1-standard-16). OS: CentOS 7; Workload: <https://github.com/gatk-workflows/gatk4-genome-processing-pipeline>.

N2: Testing completed by Intel on 9/11/2020. Region: Google Cloud us-central1 region. N2-standard-2 (2 vCPUs); N2-standard-16 (16 vCPUs). 180+ instances. 2nd Gen Intel Xeon Scalable processor (Cascade Lake); memory/instance: 7.5 GB (N1-standard-2); 60 GB (N1-standard-16). Network-attached storage; network bandwidth/instance: 10 Gbps (N1-standard-2); 32 Gbps (N1-standard-16). OS: CentOS 7; Workload: <https://github.com/gatk-workflows/gatk4-genome-processing-pipeline>.

³<https://allofus.nih.gov/>.

Performance varies by use, configuration and other factors. Learn more at www.Intel.com/PerformanceIndex.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Your costs and results may vary.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy. Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others. 102021/RJM/J RL/PDF Please Recycle xxxxx-001US

