

A New Sequence for Bioinformatics HPC

The UK National Health Service is collaborating with university partners to transform public health and personalized medicine with the power of high performance computing.



GIG
CYMRU
NHS
WALES

Healthcare



Business needs

Bioscience teams at NHS and the UK's Cloud Infrastructure for Microbial Bioinformatics project need high performance computing and storage systems to enable next-generation genome sequencing technologies.

Solutions at a glance

- Dell EMC PowerEdge C6525 servers
- PowerSwitch S3048-ON
- Mellanox® IB-SB7790
- Red Hat® Ceph Storage

Business results

- Accelerating genomic sequencing
- Fighting infectious diseases
- Improving public health
- Enabling personalized medicine

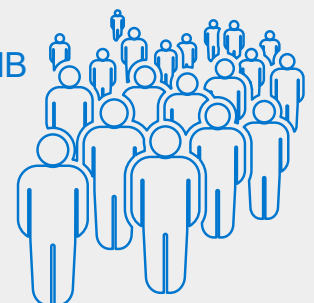
Over the past 12 months, Dr. Connor and his colleagues have sequenced

8,000 - 9,000
genomes



The HPC resources operated by MRC CLIMB are leveraged by

1,000
users



Personalizing treatment options for infectious diseases

High performance computing (HPC) has long been a staple across the broad sets of workloads in bioinformatics. However, as data volumes grow and clinicians need ultra-fast results, pure compute performance becomes a secondary concern over stability, security, storage flexibility, application portability and management.

Going from raw sequence data (millions of lines of short DNA fragments) to a final report a clinician can use for treatment requires several steps along the workflow path. A large sequencing-based effort like the Cloud Infrastructure for Microbial Bioinformatics (CLIMB) project in the UK, which gathers and analyzes gene sequences from infectious microbes to provide transmission and personalized treatment data for many of the UK's National Health Service (NHS) programs, highlights many of these challenges — and solutions. For the team behind this effort, the infrastructure considerations are as wide-ranging as they are mission-critical.

The value of a national program to monitor and provide personalized treatment options for infectious diseases is incredible, according to Dr. Thomas Connor, principal investigator for the CLIMB program's Cardiff division. The difficulty is that the computational and storage environments have demands that go beyond what some HPC research shops are set to deliver, requiring an innovative rethink of how to use HPC, cloud, scalable storage and containers.

In short, this innovative project required extraordinary creativity to navigate genomics-specific constraints.

The CLIMB project

The CLIMB project is a collaboration among Warwick, Birmingham, Cardiff, Swansea, Bath and Leicester Universities and the Quadram Institute Bioscience. It is dedicated to developing and deploying a world-leading cyber-infrastructure for microbial bioinformatics, providing cloud-based compute, storage and analysis tools for microbiologists

MRC

Cloud Infrastructure
for Microbial
Bioinformatics

across the UK, accompanied by a wide range of bioinformatics training activities.

Architecting for mission-critical bioinformatics

Genomics-based disease analysis, control and targeted treatments are separate steps in the CLIMB workflow and all deliver different elements to clinicians and health monitoring organizations nationally and globally. It is no surprise that this means separate clusters with differing capabilities that can scale securely and reliably. That is a far larger task than it may seem, Connor says, pointing to the various ways his teams have had to re-think their choices for connecting to and using storage, virtualization, containerization and overall compute allocations.

What is most surprising about this sophisticated, bifurcated need for infrastructure is that a surprising number of the steps have been automated through the use of advanced resource and workflow management tools.

For instance, at Connor's lab, a 160-core HPC cluster sits next to the sequencing system to do some initial results generation. This is necessary because for most laboratory environments, on-site networks are limited and cannot handle the massive data transfer. This cluster handles some of the pre-processing and generates initial results before being pushed to other sites that can do more sophisticated crunching.

At the other end of this first analysis run, close to the sequencing instrument, the data is shuttled to an OpenStack-based cluster that can provide other testing environments, and has backup and resilience systems in place. This way, if a lab goes down, it is still possible to get patient data without interruption, and do fast analysis. The most critical piece of this infrastructure is not compute, as much as fast, efficient storage. To meet these unique storage requirements, Connor and his colleagues worked with Dell Technologies to hook the cluster into 7 petabytes of Ceph-based storage with a Red Hat management layer.

CLIMB has become an essential national capability for microbiologists in the UK, serving more than 900 users and over 300 research groups scattered across at least 85 research institutions — from Edinburgh to Exeter, from Belfast to Norwich — spanning universities and government agencies in the UK.

“Ceph underpins everything we do at our production sites and storage is critically important,” Connor explains. “Above all, it ensures that everything we do is reproducible.”

While reproducibility is important in many academic areas, genomics workloads like those under the CLIMB project require bulletproof strategies for this. As part of their assurances to be an accredited facility, all results that are fed into NHS and World Health Organization (WHO) programs must be able to be replicated exactly, delivering the same results consistently. Getting this right goes far beyond traditional replication strategies and storage management techniques.

Connor and team, with the creative technical expertise of their Dell Technologies HPC partners, decided they could go one step further by adding virtualization specifically for reproducibility with Singularity containers alongside a workflow management package, NextFlow, to create locked-down pipelines and processes so every aspect of the workflow remained consistent. On the ground, this made it possible for teams to capture substantial compute resources as needed outside the 160-core system in the lab and easily tap into another 2,000-core machine at the main Cardiff University data center (the “Hawk” supercomputer) and another 1,000 cores they could tap into via OpenStack for an on-premises cloud.

This creative approach to using traditional HPC alongside more enterprise-oriented technologies has given Connor and his team distinct advantages in scalability, reliability and one element we have yet to discuss, security.

Securely dealing with patient data is a concern for anyone building IT infrastructure, especially when genomic and medical history information is at the core of a project. Connor points to Dell’s integration expertise across the many layers of compute, storage, network and virtualization.

“Sequence data is in the organism and isn’t patient identifiable but broadly speaking, there is strict control of patient identification information that sits within our secured systems here and inside our member organizations,” he explains. “Our gene sequencer can automatically push this data into locations that are secured for other types of processing as well. All of this works to fulfill our information governance requirements and ensures we have all we need to be compliant.”

Security is just one aspect of overall integration, aside from those listed here. On top of all of these rigorous processing and storage pipelines is automation, something that is challenging in theory, but that has been integrated with the help of Dell Technologies HPC specialists working with Connor’s team.

Complex automation made simple

The team relies on extensive and often custom automation to manage their bioinformatics pipelines across three separate systems. To streamline automation overall, the pre-processing system in the lab and the two other clusters are defined by the automation tool Ansible and, more specifically, by using Ansible playbooks.

The team has Slurm and Ceph in the lab with 30-terabytes (SSD) for scratch and data analysis. The system has a Chron job that checks for completion of tasks, including sequencing from the instrument attached to a network within the cluster. Once that job confirms, the team uses NextFlow to take the sequencing files that have been generated and run through a defined workflow. From there, all processing steps and individual elements are packaged into Singularity containers with NextFlow orchestrating the relevant analysis to be run. This entire set of stages is hands-free, leaving researchers to focus on the work at hand, with processes picking up automatically at defined points. Aside from the valuable automation overall, the combination of NextFlow and Singularity containers is also critical to reproducibility, one of the key aspects of any validated healthcare workflow.

From this point, the other two systems in CLIMB’s compute arsenal are brought into play. One is an AMD EPYC-based PowerEdge cluster from Dell Technologies with Kubernetes to continue the automation streak. The other is an OpenStack install that is in its fifth year of operation. This is where much of the development work has been done over the course of the project.

“We’ve been spinning up Slurm clusters within the OpenStack environment and can test containers and processes,” Connor explains. “It is where we’ve built almost all the processes that have been translated over to our production system in the lab. That’s given us space to prototype. We use large nodes, each of those on the current system is 32 cores with 512 RAM, and it’s also running with a large Ceph install, giving us around 1.7 petabytes of Ceph for core storage of the system.”

All of this orchestration and automation was built with expert advice from the Dell Technologies HPC teams Connor has worked with for several years, beginning with the OpenStack cluster, all the way to the newest AMD EPYC-based cluster.

Rethinking HPC: Ultimate convergence in bioinformatics

“One of the important things we’ve seen in the last few years in HPC is this convergence, this move away from absolute compute performance and more toward high-throughput. We’re working on hundreds or thousands of samples each week. We’re less interested in fast processing on one, we want to process loads. That’s a different way of working, which means different engineering is required,” Connor says.

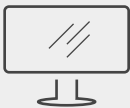
“In hospital environments, we have classic bare-metal, Slurm-based HPC, the university-based systems we use are OpenStack facing the cloud then on top of that we run Slurm in a scalable cluster on research machines that can scale with demand,” Connor adds. “We use these systems routinely. They underpin research from many groups within Public Health Wales and beyond. With these resources, we’ve been able to build clinical services for HIV and tuberculosis, influenza and C. Diff, a hospital-borne infection that can have fatal consequences. With flu, we were able to build a near-real-time service with sequencing that was pushed into the international surveillance database to inform next season’s vaccine. We were the fastest country in Europe for sure, if not the world, to do that, pushing 400 to 500 sequences last season alone.”

The steep CLIMB ahead

Connor’s teams see a vivid future on the horizon, one with AI further integrated, aided by the PowerEdge servers with NVIDIA V100 GPUs that are hosted at the Supercomputing Wales site.

Overall, Connor says they are hoping to renew their project, which finishes in 2020. There is no time like the present with the novel coronavirus at the forefront of conversations for this work to continue. He says the goal is to involve new hardware procurements, including accelerated servers to extend system capabilities and the users they support. They are also hoping for more complete integration with NHS systems to better receive and work with human pathogen data, possibly using GPUs to do detailed, secure text mining of patient information across large volumes of medical records.

The teams have done incredible work since the project began in 2016. Not just in helping public health officials and clinicians better understand and treat infectious illnesses, but in showing the way for other research centers in terms of IT infrastructure. A high-throughput oriented approach, with cloud at the core and a unique storage system, make the CLIMB project one to watch for innovation on all fronts.



[Learn more](#) about Dell EMC high performance computing solutions



[Learn how](#) to unlock the value of data with the Dell EMC solution for AI



[Share this story](#)