# THE AGE OF BIG COMPUTE

## Hybrid Cloud High Performance Computing

—

**AUTHORED BY**
**Darian Bird**
Principal Advisor, Ecosystm

# Overview

**The ability to digitise the world is increasing – from digital twins of wind turbines, simulations of autonomous vehicles, virtual screening in drug discovery, to 3D body imaging in healthcare.**

With each of these new applications comes an explosion of data and the need for high performance computing (HPC) to process it at speed. Organisations that require HPC now run more varied workloads, requiring a broader range of configurations, including different CPU architectures, high ratios of GPUs, low latency – in addition to greater memory and storage capacity. They are demanding more flexible infrastructure that can efficiently run their traditional workloads and simultaneously scale up to meet the needs of agile teams not willing to wait in job queues.

HPC-as-a-Service offerings from the hyperscalers, which provide access to massively parallel computing resources on-demand will not replace on-prem systems but will complement them. Organisations are looking to hybrid cloud HPC as a solution, for the performance and cost advantages of on-prem in addition to the elasticity and configuration benefits of the cloud. With organisations seeking to digitise more of their operations and deploy machine learning and AI applications, HPC will appeal to a wider audience. As we move into the age of Big Compute, any organisation exploring HPC should consider the benefits of a hybrid cloud system.

This whitepaper provides an overview of the benefits of HPC in the hybrid cloud, the workloads and industry use cases, and includes guidance on hybrid HPC features and deployments.

# Benefits of HPC in Hybrid Cloud

## Shift to OpEx Spending

Organisations are increasingly reluctant to approve the outlay of capital when an OpEx alternative is available. HPC in the cloud ties computing expenditure directly to use and decreases costs connected to hardware management. Meanwhile, modern on-prem offerings with consumption-based billing allow organisations to shift to OpEx spending, while still retaining much of the cost predictability associated with privately owned hardware. A hybrid system with cost management ensures resource optimisation and allows organisations more certainty around systems costs over time.

## Data Sovereignty

Many traditional HPC users are government entities or from industries that are more regulated, such as Financial Services – making data sovereignty a key buying criterion. Classified or sensitive personal data sets can only be stored or processed in multi-tenant clouds under strict conditions set out by regulators, such as APRA or the DTA in Australia. A hybrid HPC system allows organisations to maintain sensitive data in a sovereign environment and operate less-sensitive jobs in the public cloud, while continuously monitoring compliance. Moreover, by maintaining data in cloud adjacent co-location facilities, organisations can utilise cloud-based HPC without transiting through the public internet.
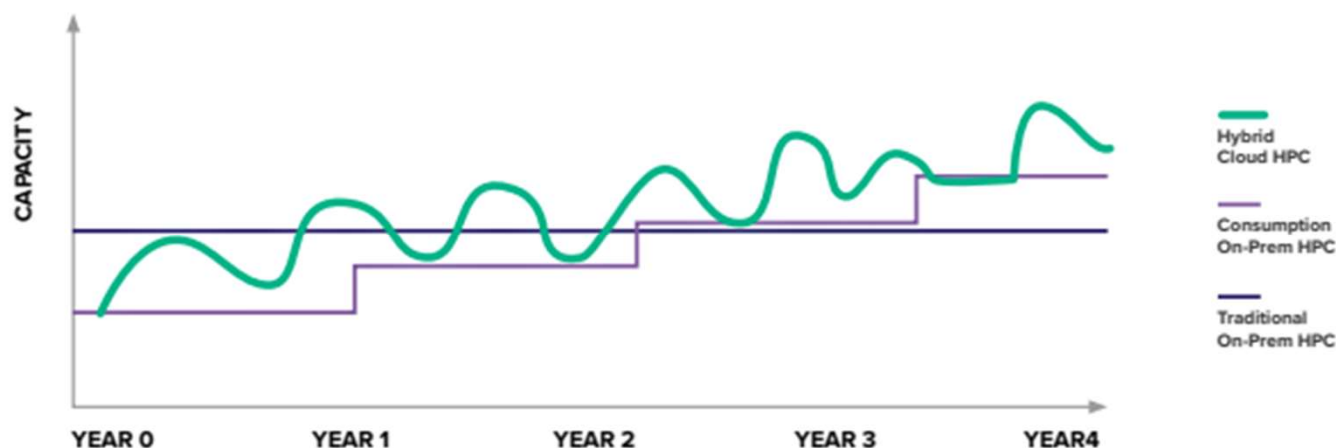
## Faster Innovation

To recoup the massive investment required for traditional HPC systems, organisations must run them continuously at full capacity. This crowded batch scheduling not only makes it difficult to add unexpected jobs but can also require waiting even for small ones at peak periods. Adding the ability to burst jobs with smaller data sets out to the cloud gives organisations greater flexibility when HPC scheduling would otherwise be a bottleneck in a workflow. An example would be a design team that uses HPC to assess a faulty product's weaknesses, and is able to run the job and iterate rather than waiting for overnight scheduling.

## Right-sized HPC Environment

Due to the high upfront cost associated with purchasing an on-prem HPC system, most organisations aim to operate them at full throttle to maximise ROI. Excess capacity is typically purchased but over the life of the machine as requirements increase, most HPC systems reach their limits. Tight scheduling can leave little headroom for unexpected jobs, particularly those that could not be anticipated years in advance during the purchase cycle. Decisions then must be made whether to invest in expansions. However, a hybrid approach allows organisations to right-size their HPC system from day one (Figure 1). An on-prem or co-located system can be sized according to current needs with additional capacity unlocked as consumption increases. If new applications demanding higher ratios of GPUs or new CPU types are required, capacity in the cloud can be scaled up and down. Advanced scheduling tools can analyse historical use patterns and simulate a range of scenarios to better understand how additional resources would impact performance and queue times.

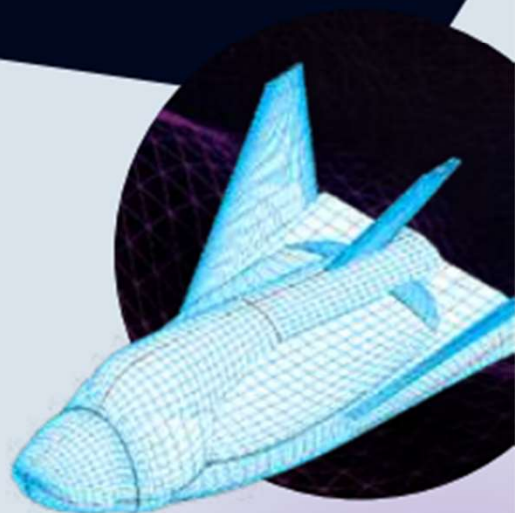### Figure 1: Hybrid HPC - Right-sized Capacity



Source: Ecosystm, 2022

# Workloads on HPC Hybrid Cloud

Organisations are increasingly deploying heterogenous workloads on HPC systems, requiring adaptable infrastructure that can be configured for a range of jobs and workflows. Some of the workloads below were traditionally only used by a small number of organisations but are now being applied across multiple industries.

## Finite Element Analysis (FEA)

One of the earliest HPC use cases, FEA generates a mesh of an object composed of smaller elements with the goal of simulating the effect of applying physical forces to it. As computational power improved, finer meshes became feasible, improving the accuracy of predictions. Moreover, implicit FEA, which is more complex and compute-intensive, has become possible allowing simulation of the impact of force over time.

## Computational Fluid Dynamics (CFD)

The analysis and solving of problems involving aerodynamics or hydrodynamics, CFD is a key tool for engineering in the Aerospace, Automotive, Construction, and Energy sectors. First, an object's geometry is derived then discretised to generate a mesh. During the solve phase, Navier-Stokes equations are typically employed to describe the effects of velocity, temperature, pressure, and density of moving fluids on the object. Finally, visualisations are analysed and refinements made to the object design.

## Monte Carlo Simulation

With thousands or tens of thousands of calculations independently conducted, the massively parallel processing capability of HPC is well-suited to speeding up Monte Carlo experiments – a mathematical method to determine probable outcomes under uncertain conditions by sampling from a probability distribution for each input. Common use cases include evaluating investment strategies, assessing network performance under various conditions, or refining weather forecasts.

## 3D Video Rendering

The need for 3D video rendering has moved beyond the Media and Gaming industries, with Real Estate, Healthcare, Manufacturing, and Geospatial Mapping all making greater use of it. HPC accelerates rendering, which requires the performance of multiple tasks simultaneously, accounting for factors such as light, texture, and perspective. A hybrid HPC system allows users to specify an optimal configuration, e.g. provisioning additional memory for scenes requiring motion blur.

# Industry Use Cases of HPC Hybrid Cloud

## Financial Services

The most significant use case for HPC in the Financial Services sector is in risk assessment. Accurately calculating risk not only reduces the likelihood of losses but also ensures compliance in an increasingly regulated industry, with Basel IV set to becoming functional in 2023. While earlier HPC systems could calculate market risk – a worst-case scenario of a portfolio's potential performance – enhancements have been needed to better understand credit and liquidity risk. To improve predictions, Monte Carlo simulations that consider multiple dependent conditions have required greater computing capacity running in parallel.

## Life Sciences

Drug development is an incredibly costly and lengthy process with a likelihood of approval from phase I studies below 10% over the last decade. Computer-aided drug design (CADD) tools, such as molecular dynamics simulation and virtual screening have only become possible with the power of HPC. Virtual screening assesses the potential of millions of molecules using the method of molecular docking to identify more diverse and active candidates. This allows researchers to rapidly move to the scoring, clustering, and selection phases. Advancements in machine learning have also created the ability to differentiate between active compounds and decoys with similar structures.

# Oil & Gas

Seismic processing and reservoir simulation are two of the most common workloads in the Oil & Gas industry to simulate wells before more expensive physical exploration takes place. Seismic imaging involves processing surface measurements to identify appropriate sites for drilling, which requires significant data storage capacity. Reservoir models are built by deploying sensors during drilling to create a 3D simulation of the well, estimating volume, porosity, saturation, and shear velocity to improve strategies for extraction.

# Manufacturing

Computer Aided Engineering (CAE) has greatly benefited from HPC by allowing manufacturers to generate high-fidelity digital twins and quickly iterate to refine designs. Aerospace, Automotive, and Industrial Manufacturing sectors use a range of HPC workloads, including FEA, CFD, and multibody dynamics (MBD) to engineer stronger, lighter, more efficient products with fewer physical prototypes. Use cases include virtual crash tests for vehicles, simulating wear and tear on moving parts, modelling airflow over wings and turbines, and studying combustion in engines.
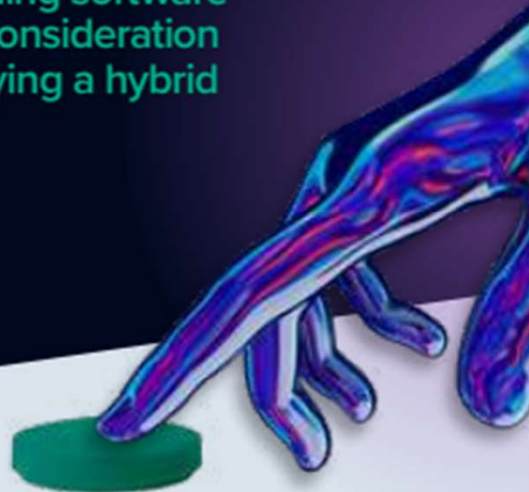
# Education

Researchers at universities have been important users of HPC to advance the knowledge of humanity. While most other industries have one or two primary HPC workloads, in Education applications vary according to the focus areas of the researchers. Traditionally, CFD, CAE, and computational chemistry were common – there has since been a growing interest in autonomous vehicle simulations, atmospheric modelling, astronomy, and genomic sequencing. Universities require an agile infrastructure with cluster management tools to pull together a broad range of configurations and orchestrate a diversity of workloads.

# HPC Hybrid Cloud Features to Consider

HPC systems are built from compute, storage, and networking components with cluster management and scheduling software to efficiently process jobs and workflows. Special consideration needs to be made for several features when deploying a hybrid HPC system.

## Chip-type Variation

One of the key technological benefits of hybrid HPC is access to a range of hardware options and configurations in the cloud that can supplement on-prem systems. In addition to CPUs by Intel and AMD, on-prem HPC vendors are increasingly incorporating ARM-based chips for greater efficiency. Since acquiring Annapurna Labs, AWS has made ARM-based Graviton chips available for HPC workloads in the cloud. Although the ratio of GPUs to CPUs in on-prem systems continues to grow, some GPU-intensive workloads, such as machine learning can benefit from additional resources delivered from the cloud. Organisations experimenting with FPGAs for application acceleration can evaluate both Intel and Xilinx hardware across different hyperscalers. Finally, Microsoft offers Cray in Azure with Cray XC or CS machines with attached ClusterStor delivered as a service for organisations looking for on-demand supercomputing.

## Batch Management and Scheduling

Traditional batch management and scheduling tools for HPC were focused on achieving as close to full utilisation of on-prem resources as possible. With a hybrid cloud system, capacity is only limited by budgets, which means management tools have become more sophisticated to most efficiently allocate workloads to resources and avoid runaway instances. Schedulers should balance total cost with queue times when directing jobs to available capacity on-prem, in reserved instances, cheaper spot instances, and if necessary, more expensive on-demand instances. Licence allocation tools allow schedulers to determine workload order more efficiently not only according to compute requirements but also the availability of software licences. One of the most popular schedulers is Slurm, an open-source, customisable workload manager.

## HPC-specific Containers

Although many organisations are already running containerised or even serverless applications, HPC workloads have been slower to shift. It is important to consider which HPC applications will eventually benefit from the portability associated with running in containers across multiple cloud environments. For DevOps teams, it may be necessary to develop new skills to take advantage of alternative container runtimes, such as Singularity, which have been designed to overcome the limitations of Docker in HPC environments. Singularity integrates natively with popular job schedulers (e.g. Slurm), supports GPUs out-of-the-box, and provides additional security by executing images without running a daemon as root.

# Recommendations

## Calculate the carbon footprint of your HPC system

Between energy-intensive compute processes and cooling requirements, HPC systems can become a significant consumer of electricity. Organisations should quantify the impact of powering their hybrid HPC systems, both in terms of the environmental and financial cost. Discuss with prospective cloud and co-location providers about renewable power purchasing agreements that they have in place. Data centre owners are also increasingly placing solar panels on the roofs of their own and neighbouring buildings to reduce emissions. For locations in which renewable energy is insufficient, carbon offset programs are becoming increasingly available. Power management to optimise efficiency is critical both on-prem and in the cloud. Energy-aware workload managers ensure that CPU frequencies are scaled up and down or servers powered down according to the scheduler. Finally, right-size your on-prem HPC system and burst to cloud when required to avoid under-utilised resources.

## Incorporate cost management into your planning

While deploying a cloud-based HPC system allows organisations to shift to OpEx spending, many stakeholders are uncomfortable with the prospect of overshooting budgets or unintentionally wasting resources. Simple cost analysis tools should be used to understand in which direction costs are trending, which services and lines of business contribute to expenditure, and to generate alerts when thresholds

# Recommendations

have been exceeded. Cluster monitoring tools help to identify orphan instances that are running unnecessarily and can be shut down. Advanced scheduling tools will assist in allocating jobs at deeply discounted spot pricing or dynamically identify those that should be run on reserved instances or on-prem according to capacity. Considering the large size of data sets involved in HPC, data storage and egress costs can become significant in the cloud. Data sitting idle should be migrated to a cheaper storage tier if frequent retrieval is not anticipated. Data optimisation is necessary to reduce egress fees, which are among the most difficult of costs to predict.

## Work with an advisor to create a workload discovery and placement roadmap

As with any cloud migration, it is important to assess each HPC workload for its suitability for migration to private cloud or one of several public clouds. HPC can particularly benefit from a multicloud approach, with each hyperscaler offering a unique range of CPUs, GPUs, FPGAs, accelerators, and connectivity. Ideally, begin by shifting straightforward, non-critical workloads with few dependencies to public cloud to gain some migration experience. Workloads which result in smaller outputs or shorter jobs that are run intermittently will also be good candidates. Eventually, build up to more complex workloads that can be run in a hybrid environment, running on-prem and cloud bursting during peak periods. Consider data gravity early in the process – moving large data sets to cloud will inevitably result in new applications being built around it. Work with an experienced advisor that can develop a roadmap, using automated discovery tools and playbooks to ease the process.