# The Case for Running AI and Analytics on HPC Clusters

**Creating a converged platform to run simulation and modeling, artificial intelligence (AI), and analytics workloads in a single cluster infrastructure supports breakthrough innovation while increasing the value and utilization of resources. Servers based on Intel® architecture are the ideal foundation for that convergence. This solution brief introduces the challenges and opportunities around running AI and analytics workloads on existing high-performance computing (HPC) clusters.**

The exponential growth in size of enterprise, academic, and government data stores over the past decade has fueled the need for resources that can transform that data from nascent potential to actionable insight. Using analytics engines such as Apache Spark*, the basis of intelligent insights has become commonplace across industries. These solutions continue to become more sophisticated as time goes on, driving more value for a variety of organizations.

One application of AI (including machine and deep learning) is to make analytics more powerful yet, with neural networks trained to predict future events based on past inputs, for example. A study by Narrative Science reports that 61 percent of respondents are currently implementing AI, with predictive analytics as the most widely used type of

AI-powered solution.[1] AI continues to become more important across a range of organizations; MarketWatch* estimates growth of the global AI market at a compound annual growth rate of 36 percent through 2024.[2]

There is a tendency in many organizations to adopt analytics platforms and AI technologies as distinct entities, rather than to focus on integrating with existing systems architectures to make business processes more effective and efficient. In many cases, dedicated new clusters may be created for deploying these capabilities, as shown in Figure 1. This approach creates data silos and the need for expensive operations related to moving and staging data. The existence of multiple clusters also sets the stage for under-utilized infrastructure, particularly for AI clusters used primarily for training deep learning networks, which tends to be sporadic.
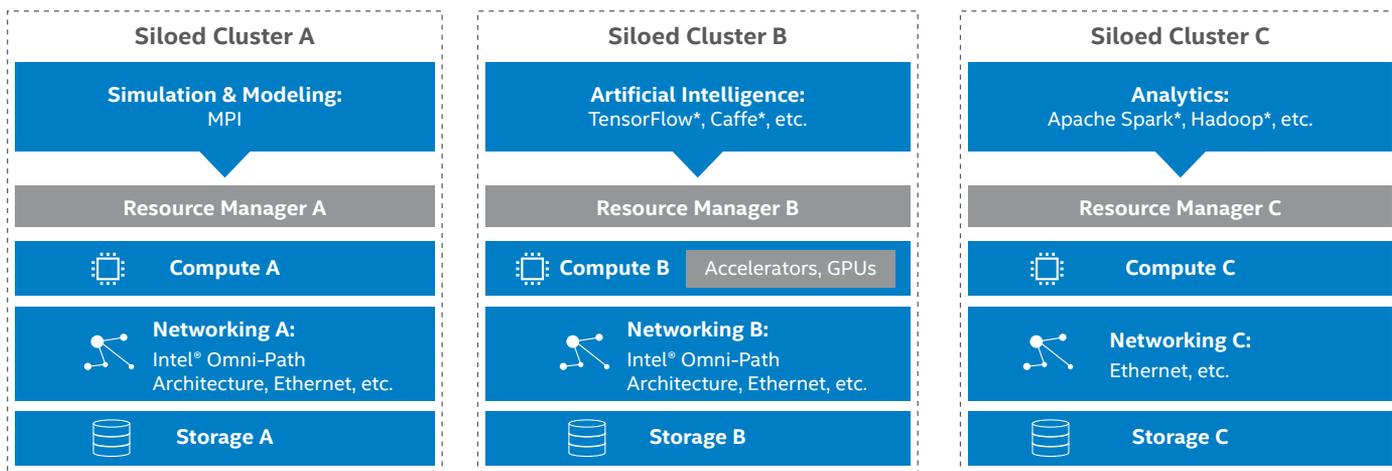


**Figure 1.** Workloads spread across separate, siloed clusters.

## Overcoming Challenges to Realize the Value of HPC

As organizations with existing HPC infrastructures adopt AI and analytics, many cite one or more of the following challenges preventing them from running all three workloads on one system:

- **Coexistent frameworks and software stacks** for simulation and modeling, analytics, and AI must be made to work together harmoniously. Unfortunately, the software stacks and frameworks for HPC, analytics, and AI are vastly different, and each workload must have its own software stack loaded on a cluster. In particular, the resource managers responsible for parceling out resources for these workloads are typically designed very differently. Intel provides a growing set of solution architectures for this purpose, as discussed in the companion brief to this one, *Supporting Simulation and Modeling, Analytics, and AI on a Common Platform.*

- **Non-HPC-oriented hardware** such as servers equipped with accelerators and GPUs draws the interest of many organizations as a potential means to drive greater performance for AI workloads in particular. In fact, the characteristics of the systems they may already own in typical HPC infrastructures (for example, robust compute cores connected with high-performance network fabrics and to high-performance shared storage) are well-suited to the needs of analytics and AI workloads. Intel has invested heavily to increase AI performance on Intel® Xeon® Scalable Processors as well as introducing new AI instructions specific to Intel® Xeon® CPUs that make them an excellent choice for AI workloads.

- **Cultural and operational separation** between HPC teams' focus on bare metal, on-premises clusters, as opposed to the cloud orientations of many AI and analytics teams, which may extend to functional approaches such as DevOps. In addition, AI and analytics teams typically use higher-level languages such as Python*, Scala*, and Java*, whereas HPC teams are more likely to be using lower-level languages such as C/C++ and Fortran.

Overcoming these challenges and the assumptions that accompany them can be an important factor in guiding an organization on the path to efficient infrastructure for all three types of workloads. The cost benefits of bringing simulation and modeling, AI, and analytics workloads together onto a single cluster infrastructure are twofold. From a capital expenditure (CAPEX) point of view, it reduces the need for new expenditures while deriving maximum benefit from existing and future cluster investments. In terms of operating expenditure (OPEX), it reduces the cost of running and maintaining the environment by simplifying the infrastructure to run on one cluster instead of multiple clusters.

These benefits are enabled through solution stacks offered by the same server original equipment manufacturers (OEMs) that IT organizations already work with as suppliers for existing HPC infrastructure. As shown in Figure 2, converged HPC stacks are poised to drive innovation and deliver value across workflows, including conventional HPC, HPC-based AI, and HPC-based analytics. Projected forward to 2021, implementations of conventional HPC are expected to continue dominating the install base, but analytics and especially AI are forecast to exhibit far higher rates of growth.[3]

## Emerging Needs for Simulation and Modeling, AI, and Analytics

As IT organizations set their architectural strategies for the next several years, many increasingly see points where AI and analytics intersect with other workloads and business processes. At the same time, the data sets that simulation and modeling jobs must handle are growing at massive rates, spurred on by current data-intensive science as well as emerging fields in cognitive computing. These trends emphasize the necessity of rich interoperability among the underlying systems, which is hampered by running simulation and modeling, AI, and analytics workloads on separate clusters. For example, consider the inefficiency of a process where simulation and modeling, data cleaning, and AI-based inference each occur on separate clusters.
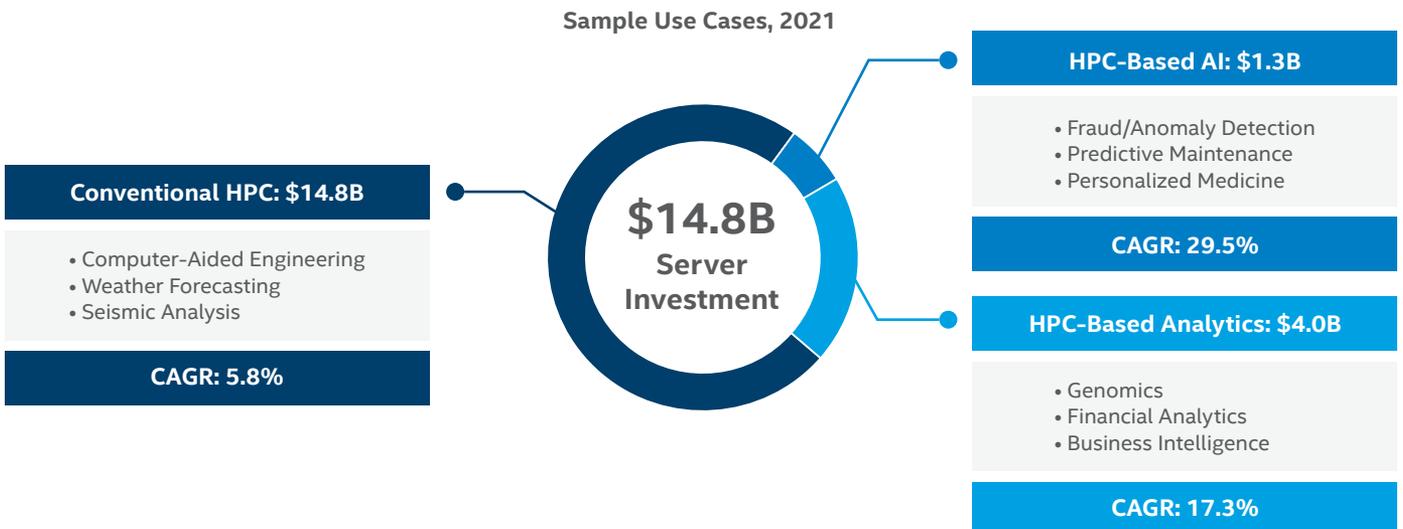
**Sample Use Cases, 2021**

**Conventional HPC: $14.8B**
- Computer-Aided Engineering
- Weather Forecasting
- Seismic Analysis

**CAGR: 5.8%**

**$14.8B Server Investment**

**HPC-Based AI: $1.3B**
- Fraud/Anomaly Detection
- Predictive Maintenance
- Personalized Medicine

**CAGR: 29.5%**

**HPC-Based Analytics: $4.0B**
- Genomics
- Financial Analytics
- Business Intelligence

**CAGR: 17.3%**

**Figure 2.** Investment levels and growth across sample HPC workloads through 2021.[3]

By running all three together on the same cluster, a number of advantages can be achieved. First, as mentioned above, the need to copy and stage data is dramatically reduced or eliminated. In addition, delay associated with transferring large amounts of data between clusters can be avoided, which is particularly important to flexibly support real-time workflows. This convergence will only become more important as emerging needs demand new capabilities within organizations.

For example, the buildout of IoT endpoints across industries will generate data flows that are orders of magnitude larger than what came before, in many cases. Sorting through those flows to identify data points of interest, patterns, and other insights is a potentially enormous task that is well-suited to AI-based analytics. Based on those outputs, opportunities may be identified for optimization of a business process or agent-based system, which could be investigated in all its permutations by simulation and modeling software. AI could be used to improve design and calibration of the actual simulation and modeling jobs.

The imperative for interoperation among the systems that cooperate across this type of overall workflow is clear. In addition, the more seamlessly these systems are integrated together, the more efficient and effective the processes can be. Successfully meeting this set of interconnected challenges paves the way for an array of use cases that span across industries and research fields, ranging from data insight, fraud detection, and robotics to individualized medicine, meteorology, and climatology.

Accordingly, many companies and institutions find themselves needing to modernize their HPC infrastructures for greater flexibility. This modernization specifically must strive to meet and bring together the needs of diverse parts of the organization, as illustrated in Figure 3. Multipurpose clusters are ideal for meeting this range of interests with maximum flexibility, efficiency, and performance. Servers based on Intel architecture are the ideal foundation for this infrastructure.

## An Open Systems Approach to Large-Scale Converged Infrastructure

Existing HPC clusters built using large-capacity, general-purpose servers based on Intel architecture can deliver flexible, cost-effective performance across workloads. In addition to eliminating the need to purchase, configure, manage, and support separate resources, this approach reduces complexity. As opposed to using systems that are designed for a specific function, jobs can be scheduled anywhere across the environment as needed, which allows for maximum utilization of resources and avoids bottlenecks associated with specific resources.

Intel works across the ecosystem—including contributing to open source projects as well as co-engineering with hardware and software providers—to make converged HPC clusters function optimally, as illustrated in Figure 4. This enablement streamlines adoption for end customers, helping them mitigate complexity and risk. It includes ensuring that resource managers and other software for simulation and modeling, AI, and analytics work together seamlessly and that frameworks are optimized for performance, scalability, stability, and security on Intel architecture. Intel also works with server providers to engineer and validate systems that deliver the best results possible.



**Business Units/ Research Departments**
Need to generate new insights and drive value by establishing innovative practices that bring together diverse data sources

**C-Suite/Board of Trustees**
Need to differentiate the organization and drive competitive advantage using the full range of available data and technology

**IT Organizations/ HPC Center Directors**
Need to efficiently deploy and manage diverse workloads and workflows as an inter-operative whole

**Users/Researchers**
Need to use new workflows and workloads that combine the capabilities of simulation and modeling, AI, and analytics

**Figure 3.** Motivation for HPC convergence across the organization.

| Simulation & Modeling: MPI | Artificial Intelligence:<br>TensorFlow*, Caffe*, etc. | Analytics:<br>Apache Spark*, Hadoop*, etc. |
|---|---|---|

**Multi-Domain Resource Managers:**
Apache Mesos*, Kubernetes*, Univa*, Slurm*, etc.

**Bare Metal Provisioning:**
xCAT*, Warewulf*, etc.

**SDI Virtualized Provisioning:**
OpenStack*, Amazon Web services* (AWS),
Microsoft Azure*, Containers

**Compute**
Intel® Xeon® Scalable processors, Intel® FPGAs, future platforms

**Storage Abstraction:** Alluxio*
Ceph*, AWS S3*, Swift*, Lustre*, POSIX*, HDFS, etc.

**Networking**
Intel® Omni-Path Architecture, Intel® Ethernet, InfiniBand*, etc.

**Figure 4.** Unified architecture across workloads.

The architecture's foundation is a common hardware and software infrastructure based on Intel® building blocks and other industry-standard components. These include a comprehensive range of popular interconnects and object stores, as well as various approaches to storage abstraction, which pull data together from all three types of workloads, creating a single data pool that is available wherever it is needed. Therefore, both the isolation of data in silos and the need to constantly move and stage large data sets among separate clusters are eliminated.

The model's unified storage architecture is based on a distributed object store that offers a number of storage abstractions for interoperability among data-access approaches. The unified fabric architecture harmonizes the behavior of various interconnect fabrics.

The architecture described here embraces both bare-metal provisioning and the virtualized provisioning approaches of software-defined infrastructure (SDI). Bare-metal approaches are the standard method used by most HPC organizations, while SDI focuses on allowing virtual resources to be spun up on demand. That dynamic approach can be instrumental in increasing agility, spawning jobs based on the constantly changing circumstances and business or research needs associated with analytics workloads.

The resource-management layer is in many ways the crux of the converged platform, providing the efficient and robust abstraction of the resources used to execute simulation and modeling, AI, and analytics workloads on the same cluster. This capability allows for the integration of functions such as resource allocation, job scheduling, and resolution of contention issues among otherwise-incompatible workloads. Domain-specific plugins at this layer allow cluster operators to tailor the operation of the environment to their needs.

This generalized stack makes it much faster and easier for organizations to explore new implementations of simulation and modeling, AI, and analytics, especially in terms of how the capabilities of all three intersect to drive new business and research value.

## Building on a Broad Intel® Technology Foundation

Rather than prescribing a rigid solution, Intel's architectural approach for simulation and modeling, AI, and analytics includes openness and flexibility as a primary design requirement. In addition to supporting hardware building blocks from the entire range of popular manufacturers, the stack draws on a vast ecosystem of software. Along with tools provided directly by Intel, many open-source and commercial third-party software packages are optimized for performance, scalability, stability, and security on Intel architecture.

Pre-validated combinations of hardware and software building blocks, tailored to specific organizational needs, are available as Intel® Select Solutions, as illustrated in Figure 5. This infrastructure helps accelerate performance while simplifying implementation and reducing risk for end customers associated with data center modernization, using systems that are available from a wide choice of popular server manufacturers.

### Intel® Architecture Building Blocks

Cluster architectures draw on a variety of Intel architecture building blocks, across the hardware stack, including the following:

- **2nd Gen Intel® Xeon® Scalable processors** are the heart of robust computing clusters that power excellent results across workloads. For example, in recent testing, it achieved an average performance improvement of up to 3.7x in HPC CPU benchmarks compared to a three-year-old system.[4] It also delivered a world-class 5.8x performance improvement on LINPACK* CPU benchmarks[5] and 1.7x better floating-point performance per core compared to competing processors.[6] These processors are also built to meet requirements of the most demanding AI inference workloads. 2nd Gen Intel Xeon Scalable processors running Intel® Deep Learning Boost have been shown to improve inference throughput by up to 25x versus the Intel® Xeon® Platinum 8180 processor.[7]

| Simulation and Modeling | | AI | | | Analytics |
|---|---|---|---|---|---|
| Intel® MPI Library | GCC and Intel® Compilers | Intel® Math Kernel Library for DNNs | Optimizations for TensorFlow* | Optimizations for MXNet* | Intel® Data Analytics Acceleration Library |
| Intel® Math Kernel Library | Optimizations for Torch* | BigDL* | Optimizations for Caffe* | Optimizations for Theano* | Optimizations for Apache* Spark |

**Tools and Optimizations Across Workloads**

| Intel® Xeon® Scalable Processors | Intel® SSDs | Intel® Omni-Path Architecture | Intel® Ethernet | Intel® Optane™ DC Persistent Memory |
|---|---|---|---|---|

**Intel® Architecture Building Blocks**

**Available as Intel® Select Solutions**

**Figure 5.** Intel® architecture hardware and software stack.

- **Intel® Optane™ DC persistent memory** offers the unprecedented combination of high capacity, affordability, and persistence. By moving and maintaining larger amounts of data closer to the processor, data-intensive HPC and AI workloads can be processed quickly and on a large scale.

- **Intel® SSDs** provide a range of storage options that enable customers to create their own balance between cost and performance. Intel® Optane™ DC SSDs deliver dramatically lower and more consistent latency, high endurance and high/balanced performance breaking through NAND SSD bottlenecks to unleash system performance.

- **Intel® Fabric products** include Intel® Omni-Path Architecture (Intel® OPA), a high-bandwidth and low-latency fabric that optimizes performance and eases deployment of HPC clusters, as well as Intel® Ethernet, an established industry leader with a broad array of options for speed, cable medium, and port count. Both are built to work smoothly in clusters with other fabrics such as InfiniBand*.

## Tools and Optimizations for Intel Architecture Across Workloads

### Simulation and Modeling-Focused Tools and Optimizations

A few key tools and optimizations that target simulation and modeling applications and workloads include the following:

- **Intel® MPI Library** is an implementation of the open-source MPICH specification designed to create, run, test, and maintain applications optimized for Intel architecture-based clusters, supporting any of multiple fabrics chosen at runtime. It provides both a runtime environment and a software development kit.

- **GCC and Intel® Compilers** enable flexible optimization capabilities for Intel architecture; the similar behaviors of both compilers allow developers to easily switch back and forth between the two. The compilers also integrate with popular toolchains in common, further aiding interoperability.

- **Intel® Math Kernel Library (Intel® MKL)** is a collection of pre-optimized math routines that are broadly applicable for technical and scientific computing. Intel maintains the functions in Intel MKL for each platform generation, enabling end customers to take advantage of hardware advances simply by relinking and recompiling their code.

- **Torch\*** is an open-source scientific computing framework with integrated support for many machine-learning algorithms. It includes a simple and robust scripting language (LuaJIT) as well as an underlying C/CUDA implementation. Intel maintains a fork of the project that is optimized for Intel Xeon processors.

### Analytics-Focused Tools and Optimizations

Following is a sample of solutions (among many others) that Intel provides or enables for improved results from analytics workloads:

- **Intel® Data Analytics Acceleration Library (Intel® DAAL)** is a collection of optimized routines to accelerate big data analytics problems. It is designed to augment the performance of applications built on popular data platforms such as Hadoop*, Apache Spark, R*, and Matlab*.

- **Apache Spark** is a particular optimization focus for Intel among analytics frameworks. Originally developed at UC Berkeley, Apache Spark is a big-data processing engine with built-in modules for streaming, SQL, machine learning, and graph processing.

*AI-Focused Tools and Optimizations*

The ecosystem for AI frameworks and other components is fast-growing, and Intel participates by helping ensure that AI implementations run best on Intel architecture, including the following:

• **Intel® MKL for Deep Neural Networks (Intel® MKL-DNN)** is a performance library for deep learning applications that provides highly vectorized and threaded building blocks for implementing deep neural networks on Intel architecture-based systems. Learn more: https://01.org/mkl-dnn.

• **Deep learning frameworks** such as BigDL*, Caffe*, MXNet*, TensorFlow*, and Theano* are optimized to accelerate AI workflows, including simplified development of applications that benefit from fast training of deep neural networks on servers and clusters based on Intel architecture. Learn more: http://intel.ai/framework-optimizations.

## Conclusion

Converging simulation and modeling, AI, and analytics platforms is a necessary evolution for enterprises and institutions that anticipate running all three types of workloads in the coming years. By supporting all three types of jobs in a single environment, many organizations will be able to reduce both CAPEX and OPEX as they embrace the full spectrum of compute capabilities that will be needed to support emerging business and research needs.

Accelerate HPC innovation:
**www.intel.com/hpc**