

TECHNICAL OVERVIEW

# INTRODUCING THE CLOUD-NATIVE SUPERCOMPUTING ARCHITECTURE



## ABSTRACT

High-performance computing (HPC) and artificial intelligence (AI) have driven supercomputers into wide commercial use as the primary data processing engines enabling research, scientific discoveries, and product development. These systems can carry complex simulations and unlock the new era of AI, where software writes software.

Supercomputing leadership means scientific and innovation leadership, which explains the investments made by many governments, research institutes, and enterprises to build faster and more powerful supercomputing platforms.

Extracting the highest possible performance from supercomputing systems while achieving efficient utilization has traditionally been incompatible with the secured, multi-tenant architecture of modern cloud computing. A cloud-native supercomputing platform provides the best of both worlds for the first time, combining peak performance and cluster efficiency with a modern zero-trust model for security isolation and multi-tenancy.

The key element enabling this architecture transition is the data processing unit (DPU). The DPU is a fully integrated data-center-on-a-chip platform that imbues each supercomputing node with two new capabilities: First, an infrastructure control plane processor that secures user access, storage access, networking, and lifecycle orchestration for the computing node, offloading the main compute processor and enabling bare-metal multi-tenancy. Second, an isolated line-rate data path with hardware acceleration that enables bare-metal performance.

## INTRODUCTION

Historically, supercomputers were designed to run a single application and were confined to a small set of well-controlled users. With AI and HPC becoming primary compute environments for wide commercial use, supercomputers now need to serve a broad population of users and to host a more diverse software ecosystem, delivering non-stop services dynamically. New supercomputers must be architected to deliver bare-metal performance in a multi-tenancy environment.

The design of a supercomputer focuses on its most important mission: maximum performance and lowest overhead. User applications leverage bare-metal performance and the fastest possible network to deliver breakthrough results and new scientific discoveries. The goal of the cloud-native supercomputer architecture is to maintain these performance characteristics while meeting cloud services requirements: least-privilege security policies and isolation, data protection, and instant, on-demand AI and HPC services.

The development of the cloud-native supercomputer architecture is based on open community development, including commercial companies, academic organizations, and government agencies. This

growing community is essential to developing the next generation of supercomputing.

The next step in that evolution is explored in the following pages—a cloud-native HPC and AI platform architecture that delivers uncompromised performance on an infrastructure platform that meets cloud services requirements.

## DPU—THE CLOUD-NATIVE SUPERCOMPUTING INFRASTRUCTURE PLATFORM

The data processing unit, or DPU, is an infrastructure platform that's architected and designed to deliver infrastructure services for supercomputing applications while maintaining their native performance. The DPU handles all provisioning and management of hardware and virtualization of services—computing, networking, storage, and security. It improves overall performance of multi-user supercomputers by optimizing the placement of applications and by optimizing network traffic and storage performance, while assuring quality of service. Moreover, it monitors cluster operations by feeding DPU-based infrastructure AI engines with telemetry data that's generated by various hardware equipment for transparent analysis and optimization. These AI engines can be used to assist in billing supercomputer users and to improve the business model of the supercomputer.

DPUs can offload infrastructure services to create a shared platform and enable security and isolation between applications based on service-level agreements in the shared environment. DPUs can also enable new NVIDIA In-Network Computing and In-Network Storage accelerations to increase overall application performance, efficiency, and scalability.

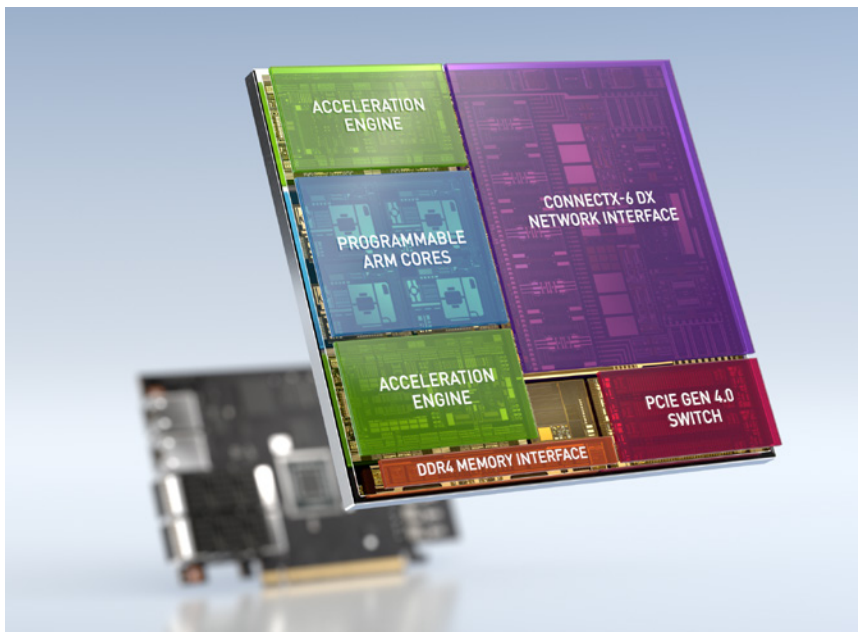


Figure 1: NVIDIA BlueField-2 DPU

DPU also support protected data computing, making it possible to use supercomputing services to process highly confidential data. The DPU architecture securely transfers data between client storage and the cloud supercomputer, executing data encryption on behalf of the user.

The NVIDIA® BlueField® DPU consists of the industry-leading NVIDIA ConnectX® network adapter, combined with an array of Arm® cores; purpose-built, high-performance-computing hardware acceleration engines with full data-center infrastructure-on-chip programmability; and a PCIe subsystem. The combination of the acceleration engines and the programmable cores enables migrating the complex infrastructure management and user isolation and protection from the host to the DPU, simplifying and eliminating overheads associated with them, as well as accelerating high-performance communication and storage frameworks.

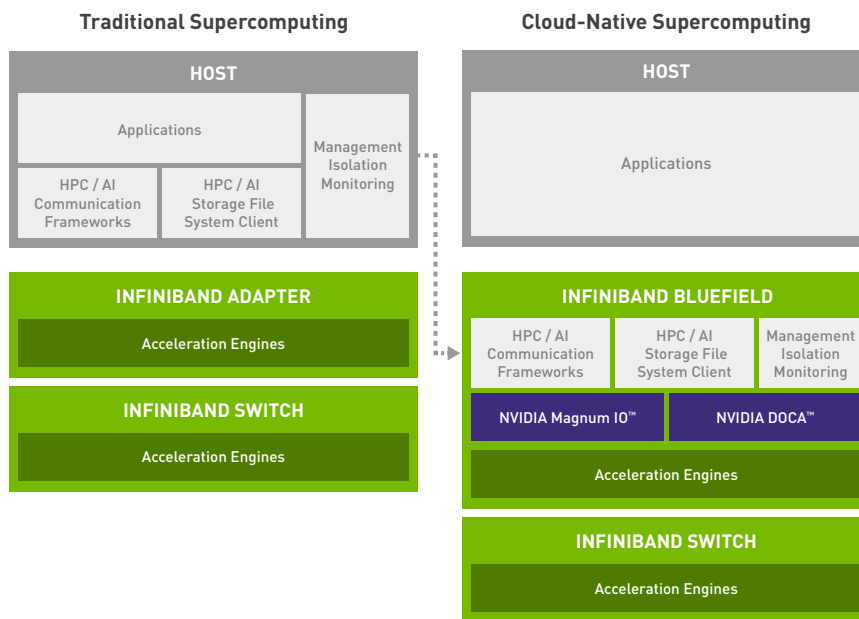


Figure 2: Cloud-native supercomputing architecture

By migrating the infrastructure management, user isolation and security, and communication and storage frameworks from the untrusted host to the trusted infrastructure control plane that the DPU is a part of, truly cloud-native supercomputing is possible for the first time. CPUs or GPUs can increase their compute availability to the applications and operate in a more synchronous way for higher overall performance and scalability. The migration of the communication and storage frameworks to the BlueField DPU also achieves a higher degree of overlapping between compute and communication, delivering the most optimal supercomputing performance and return on investment.

## MULTI-TENANT ISOLATION: TOWARD ZERO-TRUST ARCHITECTURE

The BlueField DPU enables a zero-trust supercomputing domain at the edge of every node, providing bare-metal performance with full isolation and protection in a multi-tenancy supercomputing infrastructure.

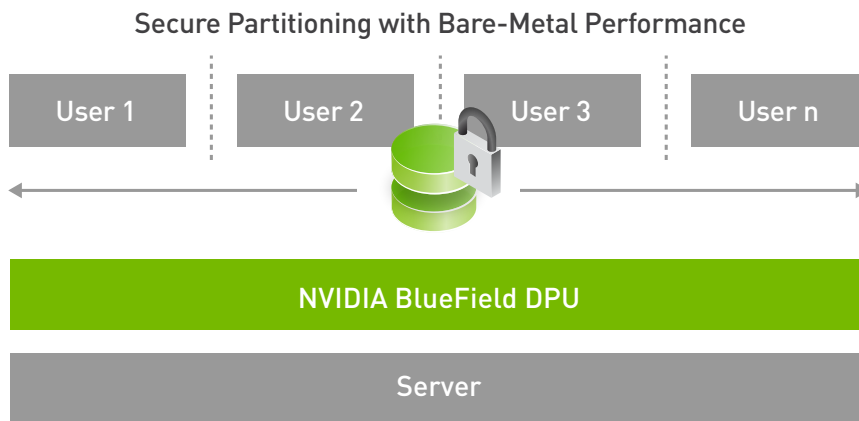


Figure 3: DPU-enabled bare-metal isolation

The BlueField DPU can host untrusted multi-node tenants and ensure that supercomputing resources used by one tenant will be handed over clean to a new tenant. As part of this process, the BlueField DPU protects the integrity of the nodes, reprovisions resources as needed, clears states left behind, provides a clean boot image for a newly scheduled tenant, and more.

## OFFLOADING HPC AND AI COMMUNICATION FRAMEWORKS

HPC and AI communication frameworks such as Unified Communication X (UCX), Unified Collective Communications (UCC), Message Passing Interface (MPI), and Symmetrical Hierarchical Memory (SHMEM) provide programming models for exchanging data between cooperating parallel processes. These libraries include point-to-point and collective communication semantics (with or without data) for synchronization, data collection, or reduction purposes. These libraries are latency and bandwidth sensitive and play a critical role in determining application performance.

A typical parallel application behavior consists of computation periods and communication periods, one after the other. This means that a new computation period cannot start before all processes have finished. Any single mode delay will delay the entire job execution by the supercomputer. In the past, communication libraries ran solely on the host CPUs, which resulted in performance bottlenecks.

Offloading the communication libraries from the host to the DPU enables parallel progress in the communication periods and in the computation periods (that is, overlapping) and reduces the negative effect of system noise. It's one of the main exascale computing strategies and is key to enabling the next generation of supercomputing architecture.

BlueField DPUs include dedicated hardware acceleration engines (for example, NVIDIA In-Network Computing engines) to accelerate parts of the communication frameworks, such as data reduction-based collective communications and tag matching. The other parts of the communication frameworks can be offloaded to the DPU Arm cores, enabling asynchronous progress of the communication semantics. One example is leveraging BlueField for MPI non-blocking, All-to-All collective communication. The MVAPICH team at Ohio State University (OSU) and the X-ScaleSolutions team have migrated this MPI collective operation into the DPU Arm cores with the OSU MVAPICH MPI and have demonstrated 100 percent overlapping of communication and computation, which is 99 percent higher than using the host CPU for this operation.

Parallel Three-Dimensional Fast Fourier Transforms (P3DFFT) is a library used for large-scale computer simulations in a wide range of fields, including studies of turbulence, climatology, astrophysics, and material science. P3DFFT is written in Fortran90 and is optimized for parallel performance. It uses MPI for interprocessor communication and greatly depends on the performance of MPI All-to-All. Leveraging the OSU MVAPICH MPI over BlueField, the OSU and X-ScaleSolutions teams have demonstrated a 1.4X performance acceleration for P3DFFT.

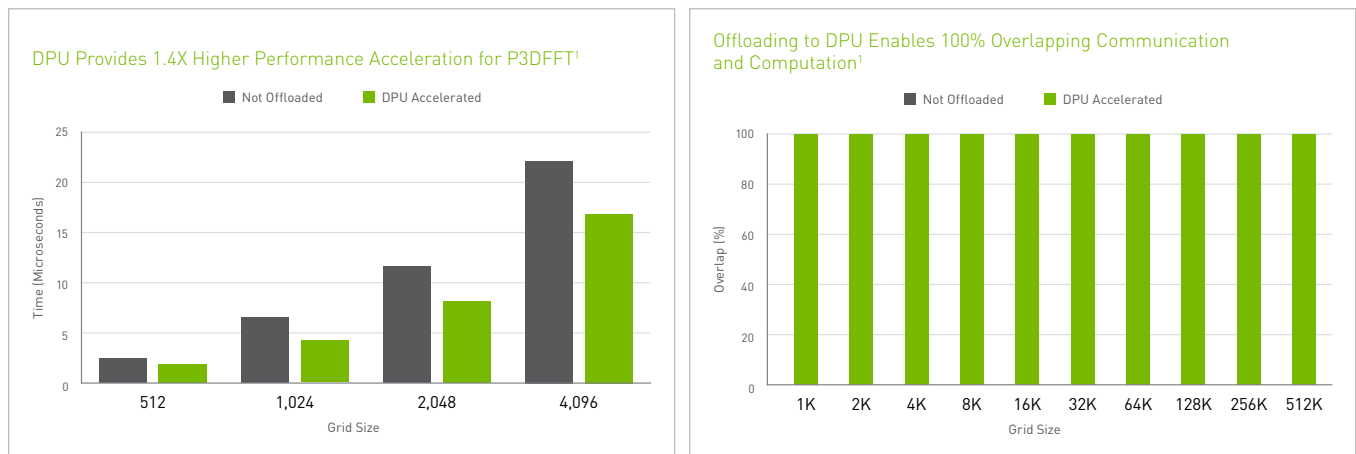


Figure 4: P3DFFT performance increase with DPU

## BLUEFIELD HPC DOCA SDK FOR CLOUD-NATIVE SUPERCOMPUTING

The NVIDIA DOCA software development kit (SDK) enables developers to rapidly create software-defined, hardware-accelerated network, storage, security, management, and AI and HPC applications and services on top of the NVIDIA BlueField DPU, leveraging industry-standard APIs. With DOCA, it's possible to program the supercomputing infrastructure of tomorrow by creating high-performance, software-defined, and cloud-native DPU-accelerated services.

DOCA provides a highly flexible environment for developing applications and services that run on DPUs while seamlessly leveraging NVIDIA In-Network Computing acceleration engines and Arm programmable engines to boost performance and scalability. In the future, embedded GPU cores will also be leveraged to execute AI algorithms on network workloads to gain increased security, performance, and more.

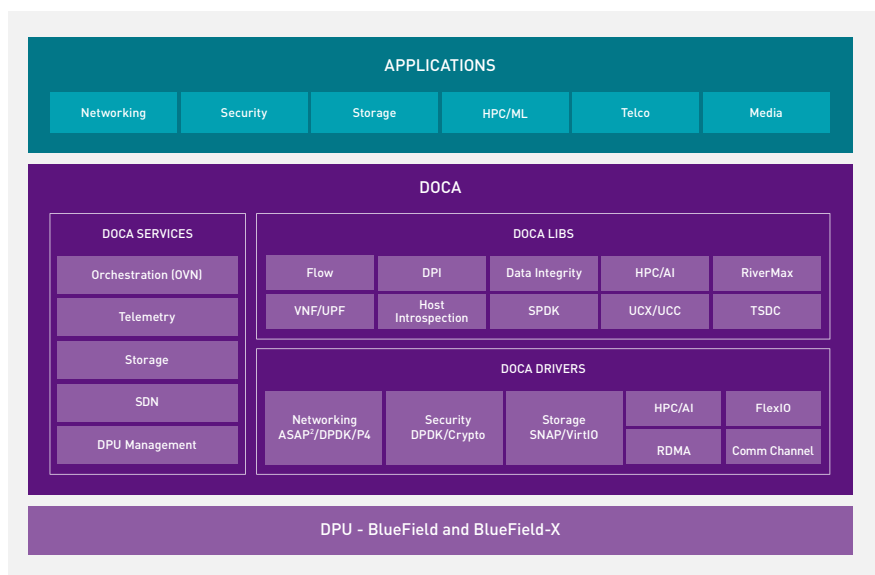


Figure 5: DOCA SDK

The DOCA SDK provides a full software stack for DPU-based supercomputers, helping to create an HPC and AI service delivery platform.

The DOCA package includes industry-standard open APIs and frameworks, such as UCX for point-to-point communications, UCC for collective communications, NVIDIA Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)<sup>™</sup> for data reductions, Open Smart Network API (OpenSNAPI), storage, security, telemetry, and management. These frameworks simplify application offload with integrated NVIDIA acceleration packages. The DOCA-based services are exposed in the compute nodes as industry-standard input/output (IO) interfaces, enabling infrastructure virtualization and isolation. The SDK supports a range of operating systems, distributions, and MPI and SHMEM libraries and includes drivers, libraries, tools, documentation, and example applications.

<sup>1</sup> The performance tests were conducted on the HPC-AI Advisory Council's Cluster Center, with the following system configuration: 32 servers with dual-socket Intel Xeon 16-core CPUs E5-2697A V4 @ 2.60GHz (total of 32 processors per node), 256GB DDR4 2400MHz RDIMMs memory, and 1TB 7.2K RPM SATA 2.5" hard drive per node. The servers were connected with NVIDIA BlueField-2 InfiniBand HDR100 DPUs and NVIDIA Quantum™ QM7800 40-port HDR 200Gb/s InfiniBand switch.

## Learn More

Learn more about **NVIDIA Cloud-Native Supercomputing Platform**

To learn more about the technology powering the platform, visit:

**[NVIDIA InfiniBand networking](#) | [NVIDIA DOCA SDK](#) | [NVIDIA BlueField DPU](#) | [NVIDIA MAGNUM IO](#)**

© 2022 NVIDIA Corporation. All rights reserved. NVIDIA, the NVIDIA logo, Base Command, BlueField, ConnectX, DOCA, DGX SuperPOD, Magnum IO, Quantum, Scalable Hierarchical Aggregation and Reduction Protocol (SHARP) are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. All other trademarks and copyrights are the property of their respective owners.

ARM, AMBA, and ARM Powered are registered trademarks of ARM Limited. Cortex, MPCore, and Mali are trademarks of ARM Limited. "ARM" is used to represent ARM Holdings plc; its operating company ARM Limited; and the regional subsidiaries ARM Inc.; ARM KK; ARM Korea Limited.; ARM Taiwan Limited; ARM France SAS; ARM Consulting (Shanghai) Co. Ltd.; ARM Germany GmbH; ARM Embedded Technologies Pvt. Ltd.; ARM Norway, AS and ARM Sweden AB. FEB22

