



NetApp HCI Overview

HCI

NetApp
October 20, 2020

This PDF was generated from https://docs.netapp.com/us-en/hci-solutions/hcvdivds_netapp_hci_overview.html on November 04, 2020. Always check docs.netapp.com for the latest.



Table of Contents

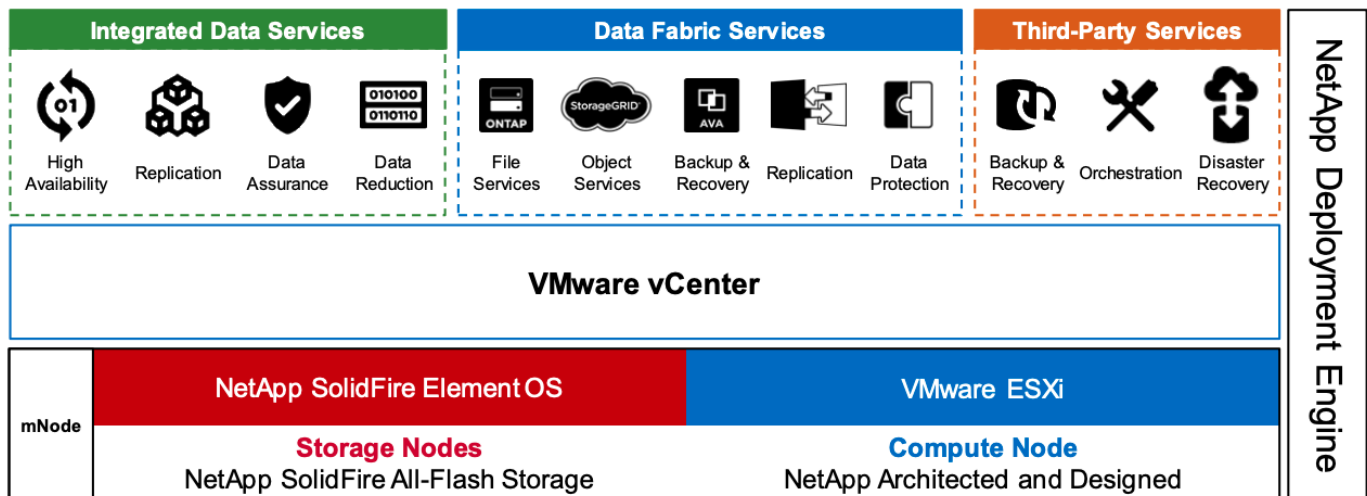
NetApp HCI Overview 1

NetApp HCI Overview

NetApp HCI is a hybrid cloud infrastructure that consists of a mix of storage nodes and compute nodes. It is available as either a two-rack unit or single-rack unit, depending on the model. The installation and configuration required to deploy VMs are automated with the NetApp Deployment Engine (NDE). Compute clusters are managed with VMware vCenter, and storage clusters are managed with the vCenter Plug-in deployed with NDE. A management VM called the mNode is deployed as part of the NDE.

NetApp HCI handles the following functions:

- Version upgrades
 - Pushing events to vCenter
 - vCenter Plug-In management
 - A VPN tunnel for support
 - The NetApp Active IQ collector
 - The extension of NetApp Cloud Services to on the premises, enabling a hybrid cloud infrastructure.
- The following figure depicts HCI components.



Storage Nodes

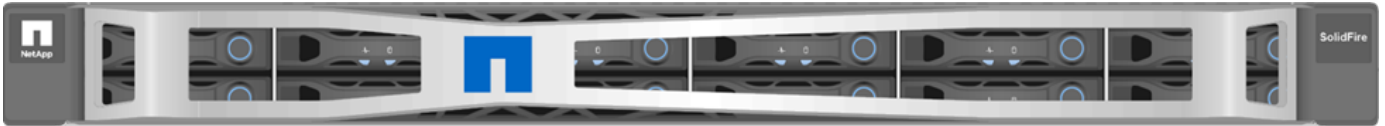
Storage nodes are available as either a half-width or full-width rack unit. A minimum of four storage nodes is required at first, and a cluster can expand to up to 40 nodes. A storage cluster can be shared across multiple compute clusters. All the storage nodes contain a cache controller to improve write performance. A single node provides either 50K or 100K IOPS at a 4K block size.

NetApp HCI storage nodes run NetApp Element software, which provides minimum, maximum, and

burst QoS limits. The storage cluster supports a mix of storage nodes, although one storage node cannot exceed one-third of total capacity.

Compute Nodes

Compute nodes are available in half-width, full-width, and two rack-unit sizes. The NetApp HCI H410C and H610C are based on scalable Intel Skylake processors. The H615C is based on second-generation scalable Intel Cascade Lake processors. There are two compute models that contain GPUs: the H610C contains two NVIDIA M10 cards and the H615C contains three NVIDIA T4 cards.









The NVIDIA T4 has 40 RT cores that provide the computation power needed to deliver real-time ray tracing. The same server model used by designers and engineers can now also be used by artists to create photorealistic imagery that features light bouncing off surfaces just as it would in real life. This RTX-capable GPU produces real-time ray tracing performance of up to five Giga Rays per second. The NVIDIA T4, when combined with Quadro Virtual Data Center Workstation (Quadro vDWS) software, enables artists to create photorealistic designs with accurate shadows, reflections, and refractions on any device from any location.

Tensor cores enable you to run deep learning inferencing workloads. When running these workloads, an NVIDIA T4 powered with Quadro vDWS can perform up to 25 times faster than a VM driven by a CPU-only server. A NetApp H615C with three NVIDIA T4 cards in one rack unit is an ideal solution for graphics and compute-intensive workloads.

The following figure lists NVIDIA GPU cards and compares their features.

NVIDIA GPUs Recommended for Virtualization

| | V100S | RTX 8000 | RTX 6000 | Available on NetApp HCI H615C T4 | Available on NetApp HCI H610C M10 | P6 |
|--------------------------------|---|---|--|--|---|---|
| |  |  |  |  |  |  |
| GPU | 1 NVIDIA Volta | 1 NVIDIA Turing | 1 NVIDIA Turing | 1 NVIDIA Turing | 4 NVIDIA Maxwell | 1 NVIDIA Pascal |
| CUDA Cores | 5,120 | 4,608 | 4,608 | 2,560 | 2,560 (640 per GPU) | 2,048 |
| Tensor Cores | 640 | 576 | 576 | 320 | — | — |
| RT Cores | — | 72 | 72 | 40 | — | — |
| Guaranteed QoS (GPU Scheduler) | ✓ | ✓ | ✓ | ✓ | — | ✓ |
| Live Migration | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Multi-vGPU | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Memory Size | 32/16 GB HBM2 | 48 GB GDDR6 | 24 GB GDDR6 | 16 GB GDDR6 | 32 GB GDDR5 (8 GB per GPU) | 16 GB GDDR5 |
| vGPU Profiles | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB, 32 GB | 1 GB, 2 GB, 3 GB, 4 GB, 6 GB, 8 GB, 12 GB, 16 GB, 24 GB, 48 GB | 1 GB, 2 GB, 3 GB, 4 GB, 6 GB, 8 GB, 12 GB, 24 GB | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB | 0.5 GB, 1 GB, 2 GB, 4 GB, 8 GB | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB |
| Form Factor | PCIe 3.0 dual slot and SXM2 | PCIe 3.0 dual slot | PCIe 3.0 dual slot | PCIe 3.0 single slot | PCIe 3.0 dual slot | MXM (blade servers) |
| Power | 250 W /300 W (SXM2) | 250 W | 250 W | 70 W | 225 W | 90 W |
| Thermal | passive | passive | passive | passive | passive | bare board |
| vGPU Software Support | Quadro vDWS, GRID vPC, GRID vApps, vComputeServer | Quadro vDWS, GRID vPC, GRID vApps, vComputeServer | Quadro vDWS, GRID vPC, GRID vApps, vComputeServer | Quadro vDWS, GRID vPC, GRID vApps, vComputeServer | Quadro vDWS, GRID vPC, GRID vApps | Quadro vDWS, GRID vPC, GRID vApps, vComputeServer |
| Use Case | Ultra-high-end rendering, simulation, 3D design with Quadro vDWS; ideal upgrade path for V100 | High-end rendering, 3D design and creative workflows with Quadro vDWS | Mid-range to high-end rendering, 3D design and creative workflows with Quadro vDWS | Entry-level to highend 3D design and engineering workflows with Quadro vDWS. High-density, low power GPU acceleration for knowledge workers with NVIDIA GRID software. | Knowledge workers using modern productivity apps and Windows 10 requiring best density and total cost of ownership (TCO), multimonitor support with NVIDIA GRID vPC/vApps | For customers requiring GPUs in a blade server form factor; ideal upgrade path for M6 |

The M10 GPU remains the best TCO solution for knowledge-worker use cases. However, the T4 makes a great alternative when IT wants to standardize on a GPU that can be used across multiple use cases, such as virtual workstations, graphics performance, real-time interactive rendering, and inferencing. With the T4, IT can take advantage of the same GPU resources to run mixed workloads—for example, running VDI during the day and repurposing the resources to run compute workloads at night.

The H610C compute node is two rack units in size; the H615C is one rack unit in size and consumes less power. The H615C supports H.264 and H.265 (High Efficiency Video Coding [HEVC]) 4:4:4 encoding and decoding. It also supports a VP9 decoder, which is becoming more mainstream; even the WebM container package served by YouTube uses the VP9 codec for video.

The number of nodes in a compute cluster is dictated by VMware; currently, it is 96 with VMware vSphere 7.0 Update 1. Mixing different models of compute nodes in a cluster is supported when Enhanced vMotion Compatibility (EVC) is enabled.

[Next: NVIDIA Licensing](#)

Copyright Information

Copyright © 2020 NetApp, Inc. All rights reserved. Printed in the U.S. No part of this document covered by copyright may be reproduced in any form or by any means-graphic, electronic, or mechanical, including photocopying, recording, taping, or storage in an electronic retrieval system-without prior written permission of the copyright owner.

Software derived from copyrighted NetApp material is subject to the following license and disclaimer:

THIS SOFTWARE IS PROVIDED BY NETAPP “AS IS” AND WITHOUT ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE, WHICH ARE HEREBY DISCLAIMED. IN NO EVENT SHALL NETAPP BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

NetApp reserves the right to change any products described herein at any time, and without notice. NetApp assumes no responsibility or liability arising from the use of products described herein, except as expressly agreed to in writing by NetApp. The use or purchase of this product does not convey a license under any patent rights, trademark rights, or any other intellectual property rights of NetApp.

The product described in this manual may be protected by one or more U.S. patents, foreign patents, or pending applications.

RESTRICTED RIGHTS LEGEND: Use, duplication, or disclosure by the government is subject to restrictions as set forth in subparagraph (c)(1)(ii) of the Rights in Technical Data and Computer Software clause at DFARS 252.277-7103 (October 1988) and FAR 52-227-19 (June 1987).

Trademark Information

NETAPP, the NETAPP logo, and the marks listed at <http://www.netapp.com/TM> are trademarks of NetApp, Inc. Other company and product names may be trademarks of their respective owners.