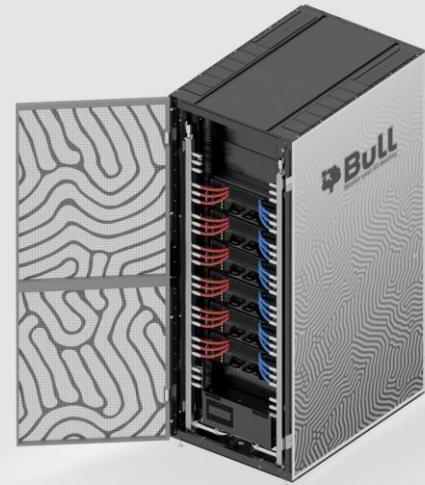


# BullSequana AI 642-H03

A high-performance AI accelerator designed for training AI models and performing inference with large models



## BullSequana AI 642-H03 overview

The BullSequana AI 642-H03 server is an AI-optimised compute node designed to empower industries, national HPC centres, government bodies, academia, and start-ups to execute accelerated, low-latency AI inference and precision fine-tuning for mid-sized language models (~100M–70B parameters) and multimodal AI workloads.

The BullSequana AI 642-H03 offers superior speed and performance, with up to 48 AMD Instinct™ MI355X GPUs per rack across six servers. The system delivers 64 TB of global GPU bandwidth, enabling high-speed communication pathways and direct memory exchange without routing through the CPU. The AMD Instinct™ MI355X GPU architecture balances GPU-CPU performance for AI workloads, supports a wide range of precision formats for training and inference including low-precision inference for mid-sized language models and offers 2.3TB HBM3e with 288GB of on-board memory per GPU.

Powered by AMD’s latest AI technology, the BullSequana AI 642-H03 combines AMD EPYC processors with high-performance GPUs, making it an ideal solution for building training clusters that require eight GPUs per server. The BullSequana AI 642-H03 server is hybrid-cooled for better energy efficiency and sustained performance across demanding workloads, featuring a manifold and a rack-mounted CDU, with up to 70% direct liquid cooling (DLC). The remaining 30% is cooled by air via fans.

## Designed for performance, built for functionality

The rack is purposely designed, customised, sized, and optimised to ensure the most effective power distribution unit (PDU), coolant distribution unit (CDU), effectively dissipating the hot air generated within compact 4U servers housed in 42U racks. For even higher energy efficiency, an optional rear-door heat exchanger can be added, targeting compute-only workloads with enhanced thermal performance.

The standardised 800 mm-wide rack allows for easier cable and accessory integration, and the in-rack CDU and manifold provide direct liquid cooling. Additionally, it features 80+ titanium-certified power supplies and PWM cooling fans.

With its fully pre-integrated, pre-tested and pre-configured design, the BullSequana AI 642-H03 server offers fast deployment, streamlined installation and reduces space and costs. This turnkey solution provides flexibility and efficiency for both rapid brownfield rollouts and greenfield deployments.

## Technical Specifications

Values are given per server; rack-level figures should be for six servers.

|                                      |  |
|--------------------------------------|--|
| <b>GPU Architecture</b>              | <ul style="list-style-type: none"> <li>Hybrid-cooled 8 GPUs AMD MI355X UBB8</li> <li>Memory: 2.3TB HBM3e / (288 GB per GPU)</li> <li>Interconnect: 8TB/s of global GPU bandwidth (1TB/s GPU-to-GPU)</li> <li>GPU TDP: 11,200W (1,400W per GPU)</li> </ul>  |
| <b>CPU Architecture</b>              | <ul style="list-style-type: none"> <li>Dual AMD EPYC™ 4th and 5th Gen processors</li> <li>12-channel DDR5 RDIMM, 3TB total size</li> </ul>   |
| <b>Storage</b>                       | <ul style="list-style-type: none"> <li>Front: 8 x 2.5" NVMe SSDs</li> <li>Internal: 2 x M.2 NVMe SSDs</li> </ul>   |
| <b>I/O Options</b>                   | <ul style="list-style-type: none"> <li>Interconnect BXL V3 / Ethernet, InfiniBand XDR and NDR</li> <li>12 x PCIe Gen5 x16</li> <li>2x 1 GbE RJ45, 2x USB 3.1, 1x VGA port and 2x 10Gb LAN ports</li> </ul>   |
| <b>Power</b>                         | <ul style="list-style-type: none"> <li>Full-server power consumption: 16 kilowatts</li> </ul>  |
| <b>Availability and RAS features</b> | <ul style="list-style-type: none"> <li>Smart Crisis Management and Protection, Dual ROM Architecture</li> <li>Hot-swap devices (PSUs, PCIe blades, NVMe drives)</li> </ul>   |
| <b>System Management</b>             | <ul style="list-style-type: none"> <li>BMC based on Aspeed AST2600</li> </ul> <p><b>Supported Operating Systems and Software</b></p> <ul style="list-style-type: none"> <li>Operating Systems: Red Hat Enterprise Linux</li> <li>Hardware Tooling: Hardware Abstraction Layer (HAL) and monitoring exporters &amp; dashboards</li> <li>Software Management Stack: <ul style="list-style-type: none"> <li>Smart Management Center xScale</li> <li>Smart Management Center</li> <li>Smart Maintenance Management Suite</li> </ul> </li> <li>Smart Performance Management Suite: <ul style="list-style-type: none"> <li>Bull OpenMPI and GPU programming environments</li> <li>Bull Slurm</li> <li>IBMS (Interconnect faBric Monitoring System)</li> <li>Singularity PRO and Enterprise</li> </ul> </li> <li>Cluster &amp; Job Monitoring <ul style="list-style-type: none"> <li>Performance and efficiency: BullSequana ARGOS</li> </ul> </li> </ul> |

As part of the BullSequana AI 600H server family, the BullSequana AI 642-H03 is meticulously assembled and rigorously tested in Bull's factory, backed by robust quality-assurance processes that ensure reliability and give customers complete confidence in the product's performance and durability.

### Software infrastructure for empowering AI teams

AI developers can utilise software infrastructure to build, train, and optimize their models effectively. By leveraging the full capabilities of the AMD ROCm 7 software stack, AI engineers can significantly enhance both usability and performance for MI355X GPUs. This enhancement leads to faster results compared to competing solutions and accelerates the model-development workflow.

AI system administrators benefit from Bull's HPC Software Suite, which provides a unified platform for running, tuning, and managing nodes and clusters at scale. It integrates orchestration, profiling, intelligent resource control, and advanced power monitoring and management to optimise performance, enhance energy efficiency, and ensure cost-effective, reliable operations. This approach aligns fully with your team's commitment to digital sobriety and sustainable infrastructure.

|                                 |  |
|---------------------------------|--|
| <b>Security Features</b>        | <ul style="list-style-type: none"> <li>Secure Boot, Secure firmware, and TPM 2.0</li> </ul>                                      |
| <b>Power Supply Unit (PSU)</b>  | <ul style="list-style-type: none"> <li>80 PLUS Titanium, Full redundancy</li> <li>AC Input: - 200-240V~/ 25A, 50-60Hz</li> </ul> |
| <b>Physical Characteristics</b> | <ul style="list-style-type: none"> <li>4U form factor</li> <li>Total Weight: 86,8 Kg</li> </ul>                                  |
| <b>Regulation and Safety</b>    | <ul style="list-style-type: none"> <li>Global Compliance: RoHS, REACH</li> <li>Certification: CE, ErP Lot 9, FCC</li> </ul>      |

### Proven track record and expertise

Bull holds a leading position as one of Europe's key supercomputer manufacturers, setting it apart in the global market. We provide end-to-end solutions – designing, delivering, installing and configuring cluster systems for leading organizations worldwide. Our factory in France plays a strategic role in this mission, enabling high-quality production that supports innovation and digital sovereignty in AI and HPC infrastructure. With decades of experience, Bull has earned the trust of industrial, academic, and HPC sectors across the globe.



Connect with us  
[bull.com](http://bull.com)



Bull is a registered trademark © Copyright 2026, Bull SAS – All rights reserved.

