

Al Infrastructure – Public Training

OUTLINE



Training Overview
In today's AI-driven era, the ability to deploy AI clusters efficiently and effectively is crucial for organizations.
This course specializes in optimizing efficiency, reliability, and scalability for deploying AI environments. It covers various infrastructure aspects, including compute platforms, storage, and the setup and maintenance of AI infrastructure.
The training focuses on key AI building blocks such as GPUs, CPUs, and BlueField networking platform, providing indepth knowledge and skills to professionals involved in deploying and managing AI infrastructure.
Training Delivery Method
Instructor-led remote training sessions via NVIDIA Teams platform.
Hands-on lab exercises focused on the data center infrastructure.
Target Audience
The course is designed for administrators, DevOps professionals, and IT-related roles who want to gain the knowledg and skills necessary to deploy and maintain AI data centers.
Training Duration

Prerequisites

Remote | 5 sessions of 5 hours

• Knowledge of core networking concepts and principles, including the TCP/IP model, Ethernet standards, basics of routing and switching, common network topologies, and IP addressing schemes.



- Hands-on experience in Linux-like systems administration, such as managing users and permissions, installing software packages, configuring network settings, and troubleshooting common issues in a Linux environment.
- Basic understanding of server hardware components and their roles in a data center environment. This includes knowledge of CPUs, memory, storage devices, and networking interfaces commonly found in servers.
- Knowledge of storage concepts and principles, including different file systems and their characteristics.
- Basic understanding of virtualization technologies, including virtual machines (VMs) and containers. You should be familiar with VM creation, management, and the role of hypervisors in virtualized environments.
- Basic understanding of artificial intelligence (AI) concepts and terminology. This may include knowledge of topics such as machine learning, deep learning, neural networks, and common AI applications.
- Before attending the course, we recommend completing the
 <u>Al Infrastructure and Operation Fundamentals</u> self-paced course. This course will provide the foundations for this training.



Training Outline

AI in the Data Center Overview

- Al Overview
- Data Center Architecture for Al Workloads

Compute Platforms for AI

- Al Compute Platforms Overview
- Hardware installation GPUs installation and Validation, Power and Cooling
- Scaling AI Compute
- Validation and Testing Cables, NCCL, HPL, NeMo
- Practice: Installing the NVIDIA GPU driver and using the nvidia-smi tool

Networking for AI

Al Data Centers Networks

BlueField Networking Platform

- BlueField Overview and Uses Cases
- BlueField Bring-up
 - Installing DOCA
 - Firmware Upgrade
 - Management via RShim
 - BlueField Interfaces Network Interfaces and OVS Bridges
- Practice: BlueField Bring-up

Storage for AI

- Storage Requirements for AI data Centers
- Storage Architecture
- Storage Configuration, Optimization and Testing
- Practice: Mounting storage and testing performance

Al Data Center Management

- Al Data Center Management Overview
- Al Infrastructure Provisioning and Management with NVIDIA Base Command Manager (BCM)
 - Overview, Components, Installation, Cluster Management Tools and Troubleshooting
- Practice: Bringing up an AI cluster with BCM



Virtualizing GPU Resources

- GPU Temporal Partitioning
- GPU Spatial Partitioning
- Virtualizing GPU resources using MIG

NVIDIA AI Software

- Using NGC Containers
- NGC CLI Installation on the Host
- NVIDIA AI Enterprise Software Suite

