

PrepAwayPDF



ONLINE TEST ENGINE

- ✓ Online Tool, Convenient, easy to study.
- ✓ Instant Online Access
- ✓ Supports All Web Browsers
- ✓ Practice Online Anytime
- ✓ Test History and Performance Review
- ✓ Supports Windows / Mac / Android / iOS, etc.

[Online Test Engine](#)



DESKTOP TEST ENGINE

- ✓ Installable Software Application
- ✓ Simulates Real Exam Environment
- ✓ Builds Exam Confidence
- ✓ Supports MS Operating System
- ✓ Two Modes For Practice
- ✓ Practice Offline Anytime

[Desktop Test Engine](#)



PDF PRACTICE Q&A'S

- ✓ Printable PDF Format
- ✓ Prepared by IT Experts
- ✓ Instant Access to Download
- ✓ Study Anywhere, Anytime
- ✓ 365 Days Free Updates
- ✓ Free PDF Demo Available

[PDF Practice Q&A's](#)



60920

Demo Downloads



59520

Successfull Cases



59062

Satisfied Clients



59146

The number of consulting

<http://www.prepawaypdf.com/>

Best Professional Test Guide Help You Pass and Provide Safe Shopping

Exam : **300-640**

Title : **Implementing Cisco Data Center AI Infrastructure**

Vendor : **Cisco**

Version : **DEMO**

NO.1 A medical company has existing third-party AI servers and NVMe storage, but they want to purchase Cisco network switches with a cloud managed feature for simplicity. Which solution best meets the requirements?

- A. Bring Your Own AI (BYO AI) HyperFabric with Cisco 6000 Series Switches and Essentials license
- B. Bring Your Own AI (BYO AI) HyperFabric with Cisco 6000 Series Switches and Premier license
- C. Nexus Dashboard with Nexus 9000 Series Switches and Cloud license
- D. HyperFabric AI with Cisco 6000 Series Switches and Premier license

Answer: B

Explanation:

Bring Your Own AI HyperFabric is the best fit when the customer already has third-party AI servers and NVMe storage but wants Cisco 6000 Series switching with simplified cloud-managed fabric operations. The Premier license aligns with the AI-focused HyperFabric use case, while Cisco describes Nexus HyperFabric as a cloud-managed AI infrastructure solution using Cisco 6000 Series switches for AI workloads.

NO.2 A company is building a network fabric for an AI training cluster that will train large language models. The cluster includes 256 GPUs distributed across 32 servers. Security requirements mandate isolation between different training projects and allowing shared access to a central storage system. Which approach provides the required security isolation and maintains optimal performance for GPU-to-GPU communication?

- A. Deploy separate physical networks for each project with dedicated switches and cross-connects to shared storage.
- B. Use VXLAN with separate VNIs per project and Group Policy Option to control granular access to shared storage.
- C. Implement port-based VLANs for each project with router ACLs controlling storage access at the distribution layer.
- D. Configure QoS-based traffic steering with separate DSCP markings per project and policy-based routing to shared storage.

Answer: B

Explanation:

VXLAN provides scalable network segmentation by assigning separate VNIs to different training projects, preserving isolation without building separate physical fabrics. Group Policy Option adds policy-based control within the VXLAN overlay, allowing controlled shared access to central storage while maintaining high-performance GPU-to-GPU communication across the AI training fabric.

NO.3 A bank is planning to deploy an AI inference solution at its branch locations. The solution must support workloads that require a balance of compute, memory, and potential GPU acceleration, and it must be suitable for installation by nontechnical onsite resources. The requirements are:

- high availability of the chassis management plane.
- support for up to 768 GB of memory for in-memory model storage and processing
- remote server launch requiring no on-site IT staff
- use of Cisco Intersight for SaaS-based infrastructure lifecycle management

Which solution meets the requirements?

- A. Cisco UCS C845A M8 Server
- B. Cisco UCS X-Series Direct
- C. Cisco Unified Edge
- D. Cisco UCS C240 Servers

Answer: C

Explanation:

Cisco Unified Edge is designed for distributed branch and edge AI deployments that need simplified remote operations through Cisco Intersight. It supports branch-friendly installation, SaaS-based lifecycle management, high-availability chassis management, up to 768 GB of memory, and balanced compute resources for inference workloads that may require GPU acceleration.

NO.4 Which set of statements describes Quantized Congestion Notification?

- A.** Priority Flow Control reacts first to mitigate congestion, and Explicit Congestion Notification acts as a fail-safe to prevent traffic drops if Priority Flow Control is insufficient.
A congestion notification packet sent toward the source by a network switch in response to receiving Explicit Congestion Notification with congestion experienced bits set.
- B.** Explicit Congestion Notification reacts first to mitigate congestion, and Priority Flow Control acts as a fail-safe to prevent traffic drops if Explicit Congestion Notification is insufficient.
A congestion notification packet sent toward the source by hosts and the network switch in response to receiving Explicit Congestion Notification with congestion experienced bits set.
- C.** Explicit Congestion Notification reacts first to mitigate congestion, and Priority Flow Control acts as a fail-safe to prevent traffic drops if Explicit Congestion Notification is insufficient.
A congestion notification packet sent toward the source by a host in response to receiving Explicit Congestion Notification with congestion experienced bits set.
- D.** Priority Flow Control reacts first to mitigate congestion, and Explicit Congestion Notification acts as a fail-safe to prevent traffic drops if Priority Flow Control is insufficient.
A congestion notification packet sent toward the source by hosts and the network switch in response to receiving Explicit Congestion Notification with congestion experienced bits set.

Answer: C

Explanation:

Quantized Congestion Notification relies on Explicit Congestion Notification as the primary congestion signal so the sender can reduce its transmission rate before packet loss occurs. Priority Flow Control acts as a safety mechanism if congestion becomes severe, and the congestion notification packet is generated by the receiving host when it receives packets marked with congestion-experienced bits.

NO.5 A financial company plans to deploy an AI training and inference infrastructure to support its customers. The team must use a solution that provides 96 GPUs, aligns with NVIDIA-compliant reference architectures, and enables cloud-based management of the entire on-premises network environment. Which Cisco solution meets the requirements?

- A. Cisco HyperFabric AI Small AI Cluster with Cisco 6000 switches and AI servers on premises
- B. Cisco HyperFabric AI Medium AI Cluster with Cisco 6000 switches and AI servers on premises
- C. Cisco Nexus Dashboard deployed in AWS or Azure with Nexus 9000 switches and AI servers on premises

D. Cisco Intersight AI fabric with Nexus 9000 switches and AI servers on premises

Answer: A

Explanation:

Cisco HyperFabric AI Small AI Cluster matches the 96-GPU requirement because the validated HyperFabric reference architecture uses 12 Cisco UCS C885A M8 servers with NVIDIA HGX for a total of 96 GPUs. Cisco Nexus HyperFabric is also cloud-managed and compliant with NVIDIA Enterprise Reference Architecture, enabling centralized management of the on-premises AI network, GPU servers, and storage.

NO.6 A server administrator must monitor the performance and utilization of GPUs in Cisco UCS servers running AI workloads. Which two metrics must be used in Cisco Intersight to accomplish these goals? (Choose two.)

A. GPU Utilization

B. Explorer

C. Monitor

D. GPU Performance

E. System CPU Utilization

Answer: AD

Explanation:

Cisco Intersight provides GPU-focused telemetry for UCS servers, including utilization and performance metrics. These measurements allow the administrator to monitor how heavily the GPUs are being used and whether the AI workload is receiving the expected GPU performance.

NO.7 A customer deploys a UCS 885a server in a data center with two power grids. Three of the 3000W power supplies in the server connect to each grid. Due to an external event, one of the grids will be shut down for maintenance. How does the server respond to this event?

A. The server immediately shuts down and stops all operations.

B. The server operates with reduced performance and caps power to GPUs at 60%.

C. The server continues to operate at full performance using redundant power supplies, but the PSU status LED glows solid amber.

D. The server initiates a graceful shut down of the operating system.

Answer: B

Explanation:

With three 3000W power supplies lost from one power grid, the UCS C885A continues operating on the remaining grid but no longer has the full power budget available for maximum GPU operation. The server reduces performance by enforcing power limits, including capping GPU power, so workloads can continue running within the available power capacity.

NO.8 A network operations team deploys HyperFabric AI and wants to monitor fabric health to proactively identify performance issues. The team must track metrics specifically relevant to AI training workloads. Which metric is most critical for identifying congestion issues that could impact AI training performance in a HyperFabric AI deployment?

A. spanning tree topology change notifications

B. BGP EVPN route count per leaf switch

- C. VXLAN tunnel packet encapsulation rate
- D. PFC pause frame transmission rate per interface

Answer: D

Explanation:

PFC pause frame transmission rate is the key congestion indicator for lossless Ethernet fabrics used by AI training workloads. A high or increasing pause frame rate shows that interfaces are experiencing congestion and applying flow control, which can directly affect GPU-to-GPU communication performance in a HyperFabric AI environment.