

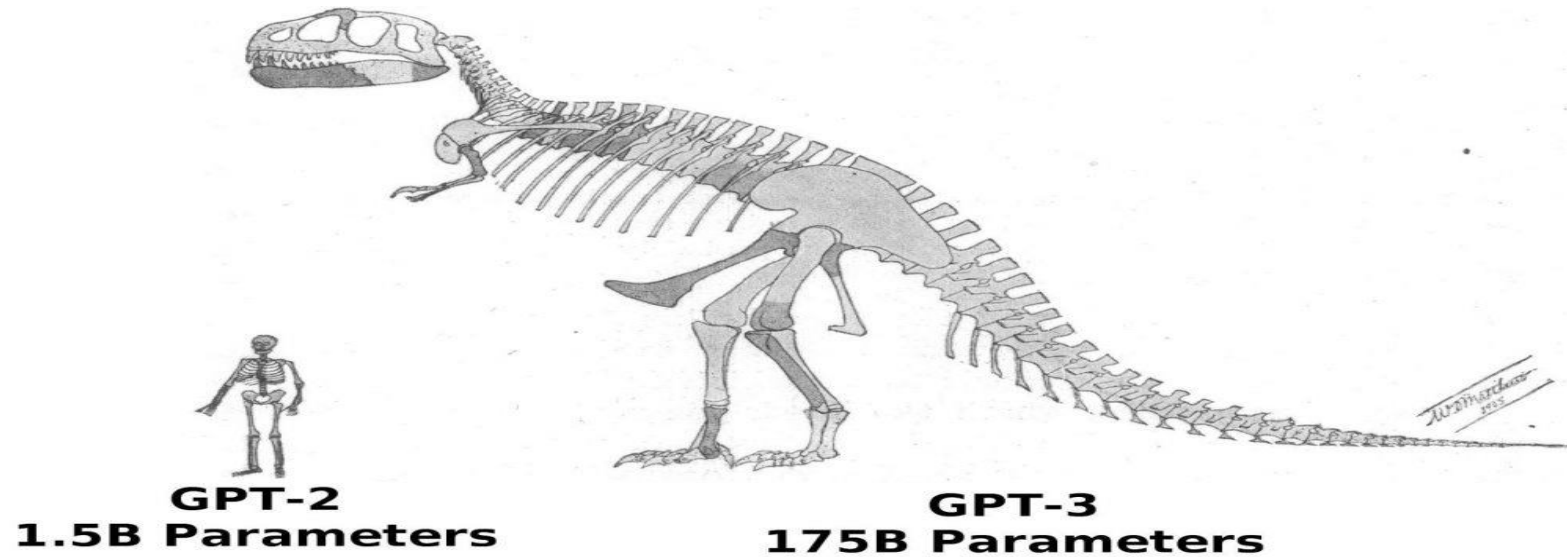


NVIDIA DGX A100 SUPERPOD 소개



BayNex

Huge data and Huge Model



# of GPUS with IB	100	200	500	1,000	2,000	5,000	10,000
Training days (V100)	1730.6	865.3	346.1	173.1	86.5	34.6	17.3
Training days (A100)	641	320.5	128.2	64.1	32	12.8	6.4

Training time projection (V100 / A100)

영상보기

<https://1drv.ms/v/s!ApvKB8j96v0lgp181efTQciST8qG8A?e=3ESD1D>

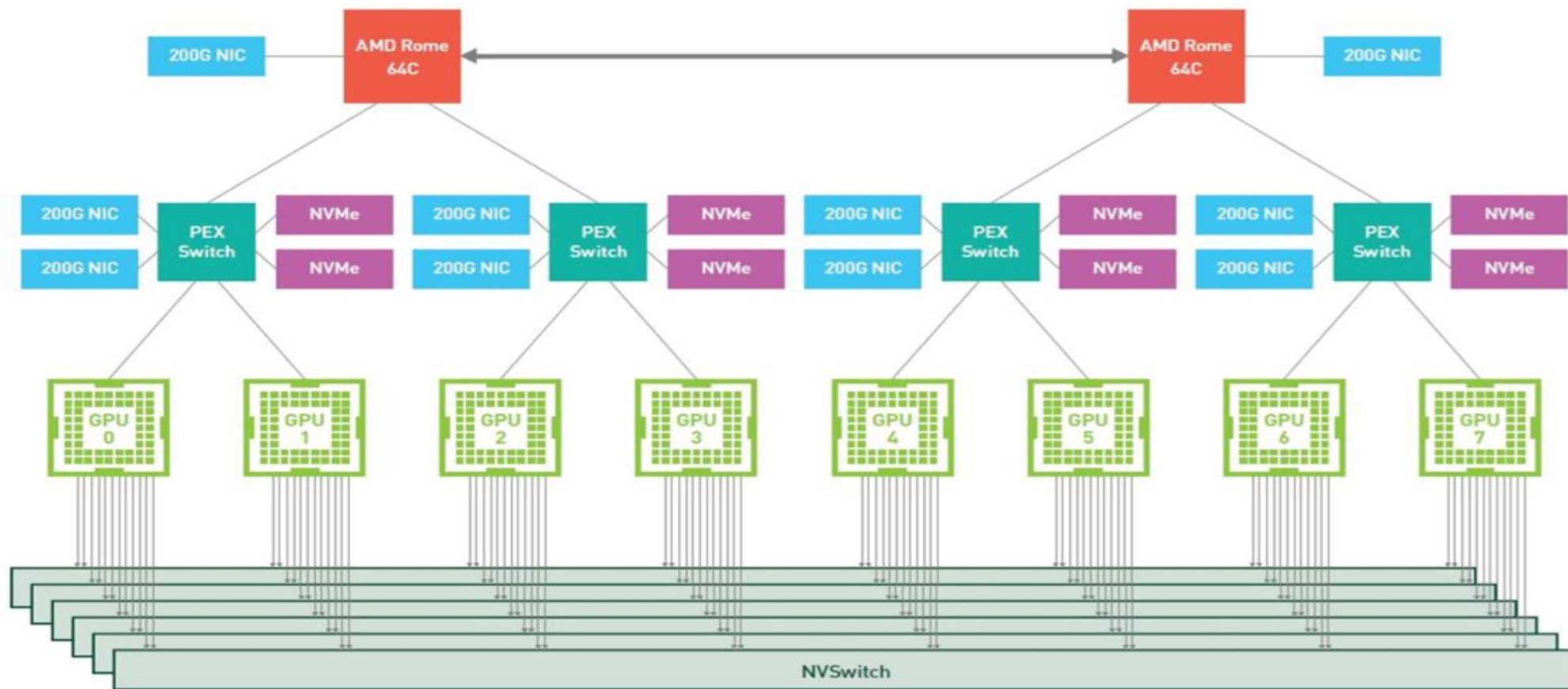


SuperPOD Architecture



 NVIDIA



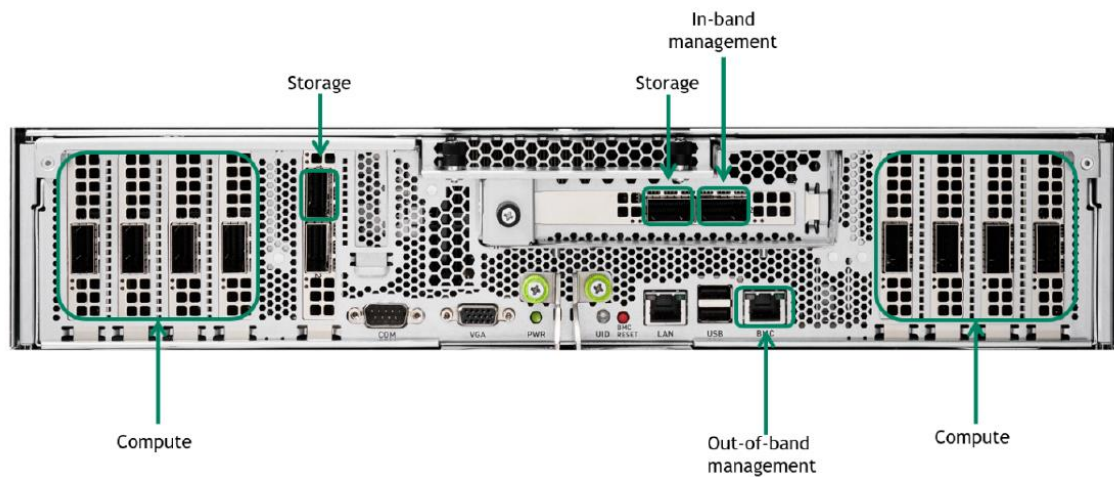


GPUs	8x NVIDIA A100
GPU Memory	320 GB total
Peak performance	5 petaFLOPS AI 10 petaOPS INT8
NVSwitches	6
System Power Usage	6.5kW max
CPU	Dual AMD Rome 7742 128 cores total, 2.25 GHz(base), 3.4GHz (max boost)
System Memory	1TB
Networking	8x Single-Port Mellanox ConnectX-6 200Gb/s HDR Infiniband (Compute Network) 1x (or 2x*) Dual-Port Mellanox ConnectX-6 200GB/s HDR Infiniband (Storage Network also used for Eth*)
Storage	OS: 2x 1.92TB M.2 NVME drives Internal Storage: 15TB (4x 3.84TB) U.2 NVME drives
Software	Ubuntu Linux OS (5.3+ kernel)
System Weight	271 lbs (123 kgs)
Packaged System Weight	315 lbs (143 kgs)
Height	6U
Operating temp range	5°C to 30°C (41°F to 86°F)

* Optional upgrades



Highest Network Throughput for Data and Clustering



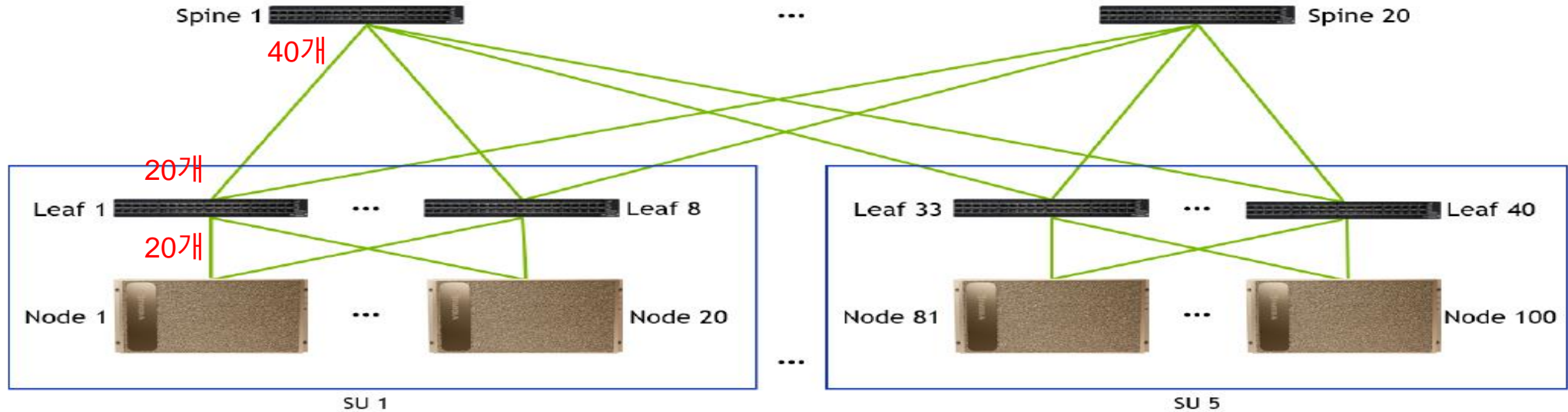
- ▶ **Compute fabric.** Connects the eight NVIDIA Mellanox ConnectX-6 HCAs from each DGX A100 through separate network planes.
- ▶ **Storage fabric.** Uses two ports, one each from two dual-port ConnectX-6 HCAs connected through the CPU.
- ▶ **In-band management.** Uses a 100 Gbps port on the DGX A100 system to connect to a dedicated Ethernet switch.
- ▶ **Out-of-band management.** Connects the baseboard management controller (BMC) port of each DGX A100 system to an additional Ethernet switch.

SU(1) = A100 20Ea

DGX A100 scalable unit(SU)



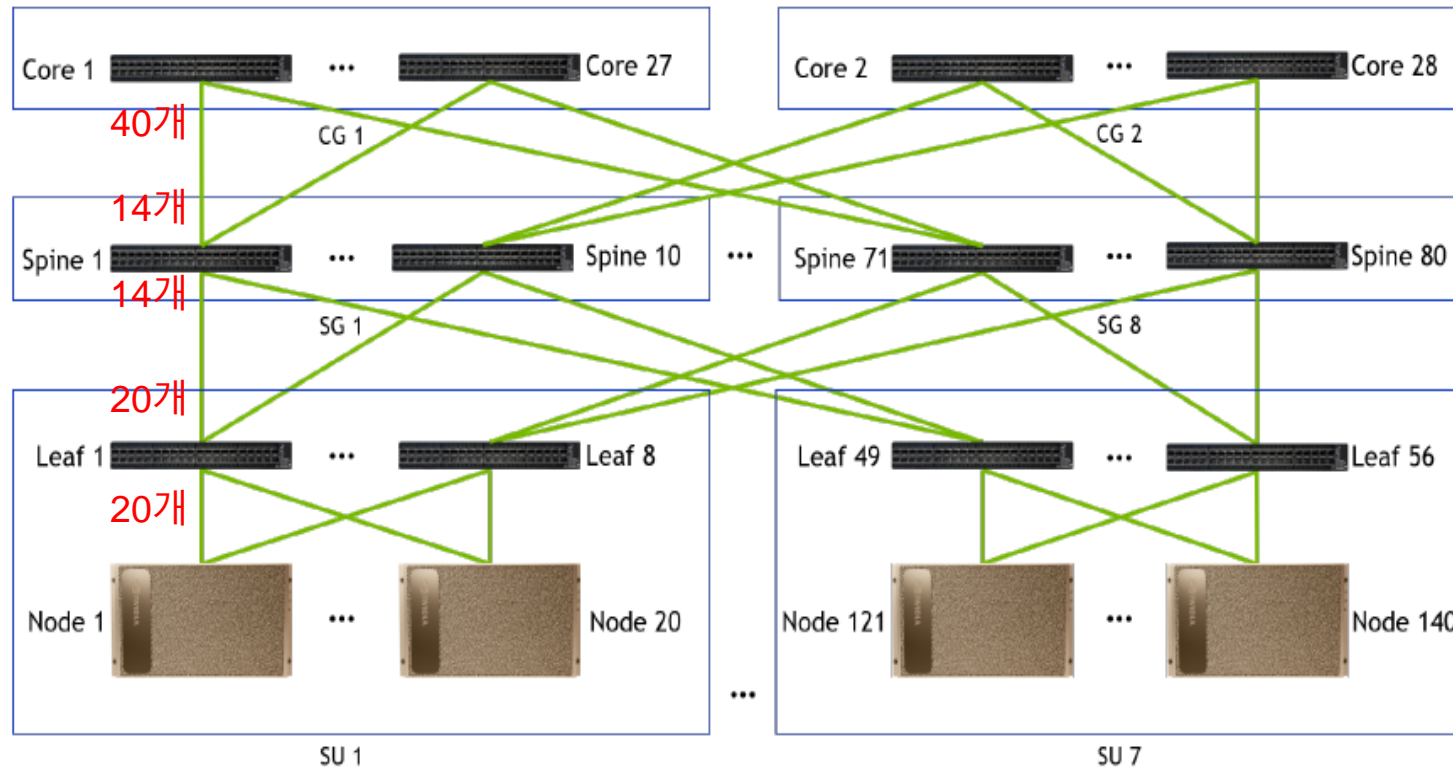
100 nodes or fewer is simpler as the third layer of switching is not required



Nodes	SUs	QM8790 Switches			Cables		
		Leaf	Spine ¹	Core ¹	Leaf	Spine	Core
10	½	8	2		80	80	
20 (Single SU)	1	8	4		160	160	
40	2	16	10		320	320	
80	4	32	20		640	640	
100	5	40	20		800	800	
120	6	48	80	24	960	960	960
140 (DGX SuperPOD)	7	56	80	28	1120	1120	1120

1. To avoid possible performance issues, ports on Spine and Core switches must only be used for inter-switch cabling.

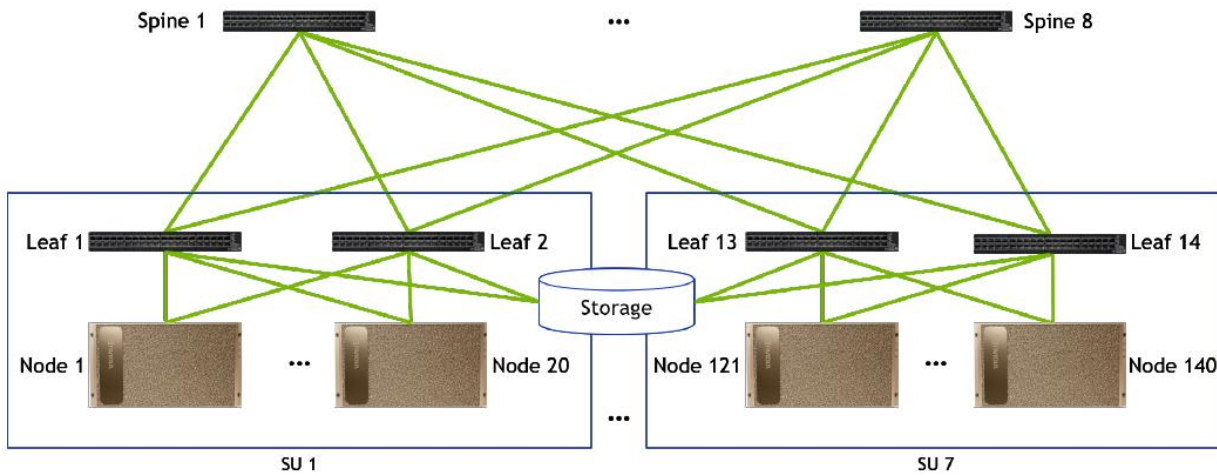
140 nodes



Nodes	SUs	QM8790 Switches			Cables		
		Leaf	Spine ¹	Core ¹	Leaf	Spine	Core
10	½	8	2		80	80	
20 (Single SU)	1	8	4		160	160	
40	2	16	10		320	320	
80	4	32	20		640	640	
100	5	40	20		800	800	
120	6	48	80	24	960	960	960
140 (DGX SuperPOD)	7	56	80	28	1120	1120	1120

1. To avoid possible performance issues, ports on Spine and Core switches must only be used for inter-switch cabling.

Storage Fabric



Nodes	SUs	Storage Ports	QM8790 Switches		Cables		
			Leaf	Spine	To-Node	To-Storage	Spine
10	½	4	2	1	20	4	16
20	1	8	2	1	40	8	32
40	2	16	4	2	80	16	64
80	4	32	8	4	160	32	96
100	5	40	10	4	200	40	160
120	6	48	12	6	240	48	192
140	7	56	14	8	280	224	56

Compute , Storage Network

. QM8790 switch



40 QSFP56 ports (50G PAM4 per lane)

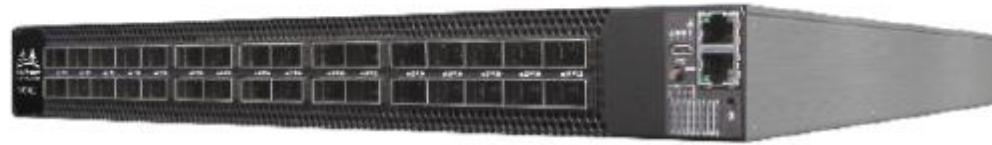
- 40 ports of HDR, 200G
- 80 ports of HDR100, 100G

Superior performance

- 90ns latency
- 390M packets per sec (64B)
- 16Tb/s aggregate bandwidth

In-Band Management Network

SN3700C switch



- ▶ Connects all the services that manage the cluster.
- ▶ Enables access to the home filesystem and storage pool.
- ▶ Provides connectivity for in-cluster services such as Slurm and Kubernetes and to other services outside of the cluster such as the NGC registry, code repositories, and data sources.

Out-of-Band Management Network

AS4610 switch



The out-of-band Ethernet network is used for system management via the BMC and provides connectivity to manage all networking equipment. Out-of-band management is critical to the operation of the cluster by providing low usage paths that ensure management traffic does not conflict with other cluster services.



DataCenter Configuration

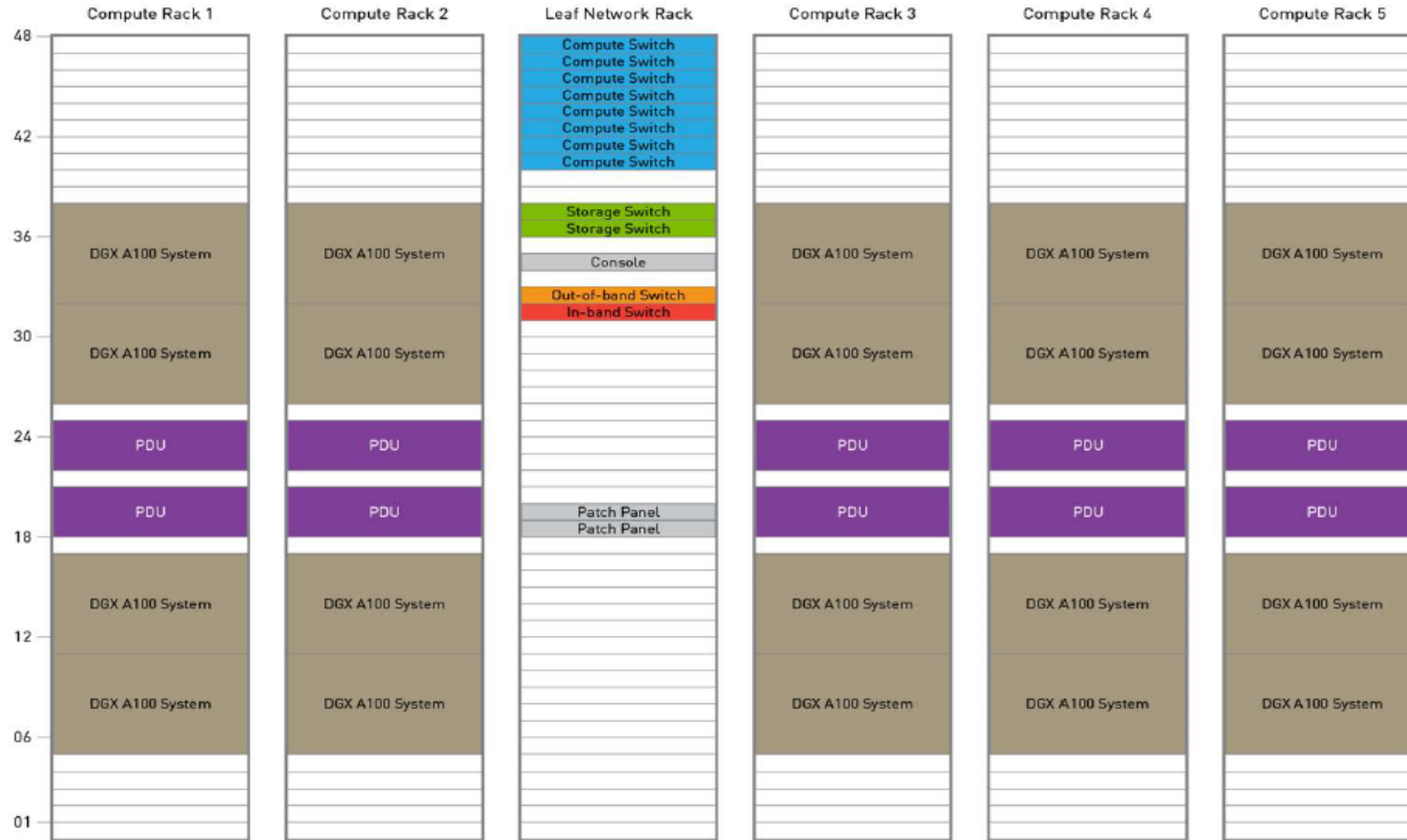
Power

Equipment	Maximum Power
DGX A100 system	6.50 kW
Management nodes	0.60 kW
QM8790 switch	0.65 kW
SN3700C switch	0.50 kW
MAS4610 switch	0.10 kW

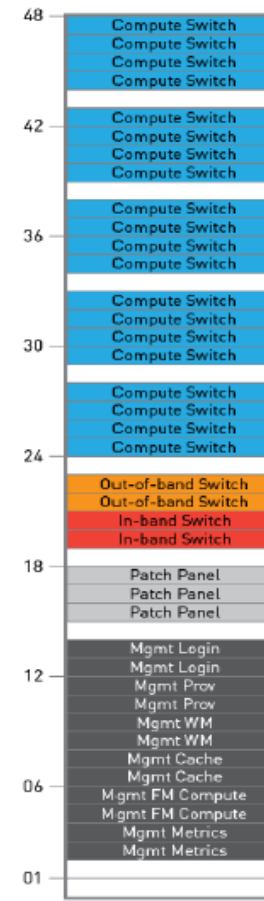
각 SU 당 137 Kw
Rack당 26 Kw

Superpod는 전체 1 Mw (스토리지 20 Kw가정)

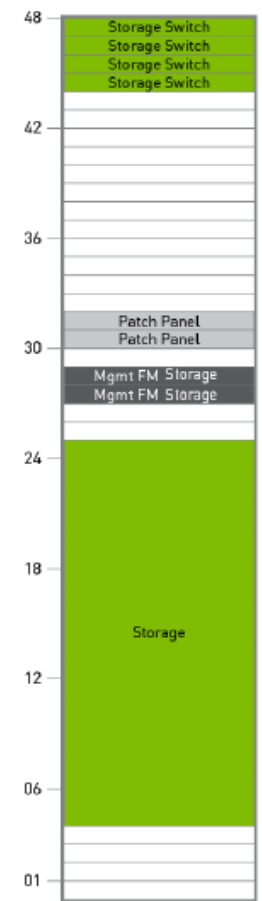
Compute: Scalable Unit (SU)



Compute Fabric and Mgmt



Storage



DGX SUPERPOD WITH DGX A100 SYSTEMS

NVIDIA DGX A100 System and Features of the DGX SuperPOD



Component	Technology	Description
Compute Nodes	NVIDIA DGX A100 System	<ul style="list-style-type: none">• 1120 DGX A100 SXM4 GPUs• 45.6 TB of HBM2 memory• 336 AI PFLOPS via Tensor Cores• 140 TB System RAM• 2.2 PB local NVMe• 600 GBps NVLink bandwidth per GPU• 4.8 TBps total NVSwitch bandwidth per node
Compute Fabric	NVIDIA Mellanox Quantum QM8790 HDR InfiniBand Smart Switch	Full fat-tree network built with eight connections per DGX A100 system
Storage Fabric		Fat-tree network with two connections per DGX A100 system
In-band Management Network	NVIDIA Mellanox SN3700C switch	One connection per DGX A100 system
Out-of-band Management Network	NVIDIA Mellanox AS4610 switch	One connection per DGX A100 system
Management Software	DeepOps DGX POD Management Software	Software tools for deployment and management of SuperPOD nodes and resource management
User Runtime Environment	NGC	Containerized DL and HPC applications, optimized for performance
	Slurm	Orchestration and scheduling of multi-GPU and multi-node jobs

감사합니다.

