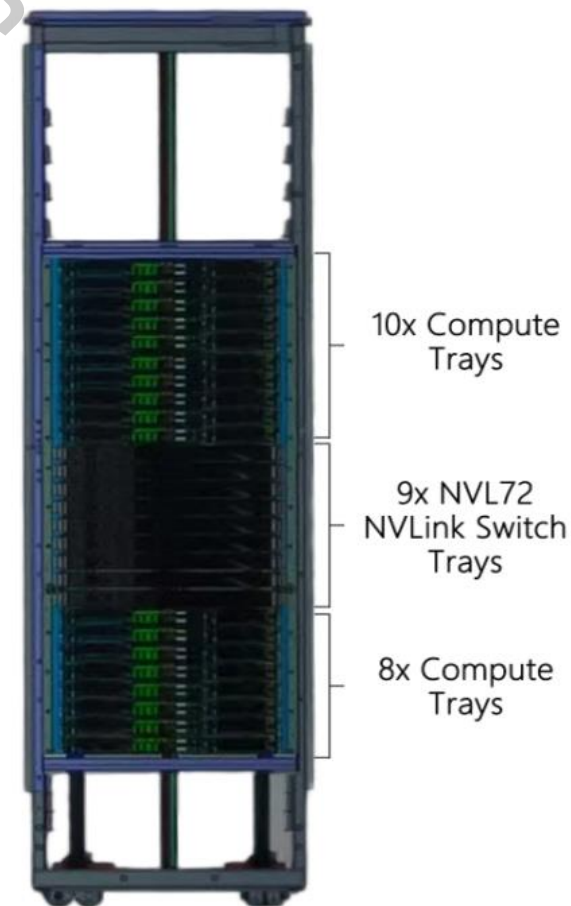
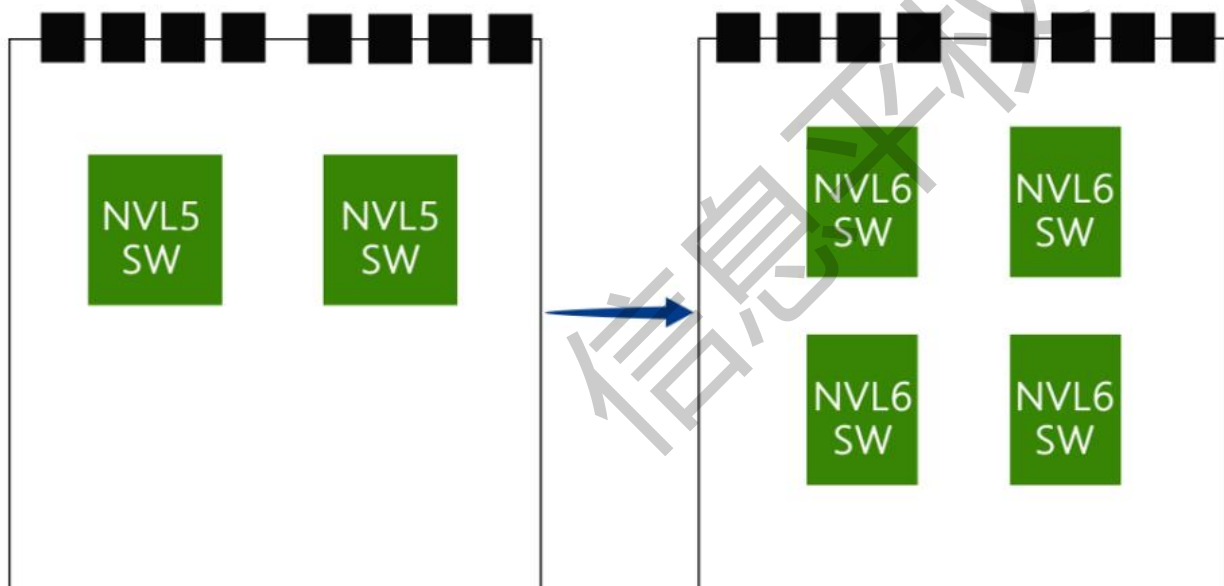
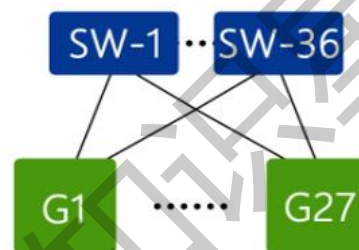


Vera Rubin NVL72

NVL Topology

- Radix of 72 means NVL72 is largest L1 domain
 - Same as Blackwell
- 1800 GB/s/dir per GPU on 36 planes
 - 2x vs. Blackwell
- 9 switch trays with 4 switch ASICs each
 - 2x switch ASICs vs. Blackwell

Note: non-scalable NVL72 pictured but scalable NVL36 also possible



Increased GPU Density

Background

- Motivation
 - Overarching goal is Perf/W and Perf/TCO
 - Part of delivering this means connecting more GPUs at higher bandwidth at minimal power/cost and highest possible reliability
 - Passive copper is lowest power, lowest cost, highest reliability, lowest risk interconnect
- Continue to drive density to enable more passive copper interconnect
- Looking at multiple vectors to achieve this
 - 100% liquid cooled-fanless, two sided service
 - Orthogonal compute/switch arrangement
 - More content on fewer PCBs-cut down cables and connectors
 - Power infrastructure-support for higher power, push AC/DC conversion out of critical volume, etc..
 - More
- Exploring concepts with 288 GPUs per rack

Vera Rubin Systems Overview

- Focus Areas for guiding datacenter infrastructure planning
 - NVL72
 - Future system concepts with dense GPU rack
- Future discussions topics
 - NVL8
 - Large NVL Air cooled configs (i.e. follow-on to GB200A NVL36)
- Requesting feedback
 - VR200 node topology
 - Denser GPU racks with higher power for long lead infrastructure planning

System Architecture

Future Roadmap

Continue to drive density to connect more GPUs at higher bandwidth with lowest power/cost and highest reliability

Rubin-Next with new high radix NVL switch + denser system design enables 4x larger NVLink domain in single large L1 and copper interconnect

- Larger effective NVL domain size

- Better binpacking/higher availability in presence of failures

High density + high radix Spectrum switch enables In-Rack L1 switch and NIC↔L1 connectivity in copper while maintaining good 2-level scale

- Lower power/cost

- Improved Reliability

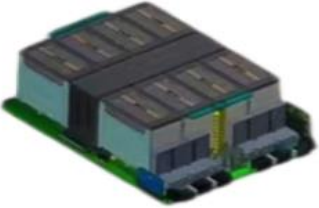
VR200 ULTRA
NVL72



288 GPU DENSE
RACK CONCEPT



System Architecture Evolution

	HGX	MGX NVL72	Future Concept
GPUs per Rack	32	72	288
Copper NVL Domain	8	72	288
Rack Power	~50 kW	~120-250 kW	~1 MW
Cooling	100% Air	~85% Liquid ~15% Air	100% Liquid
Node Size			
GPU↔CPU	8 GPUs	4 GPUs	2-4 GPUs
Interface	PCIe	C2C	C2C
PCIe Retimers	Yes	NO	NO
E/W Network NIC↔L1 Cable	Optics	Optics	Copper