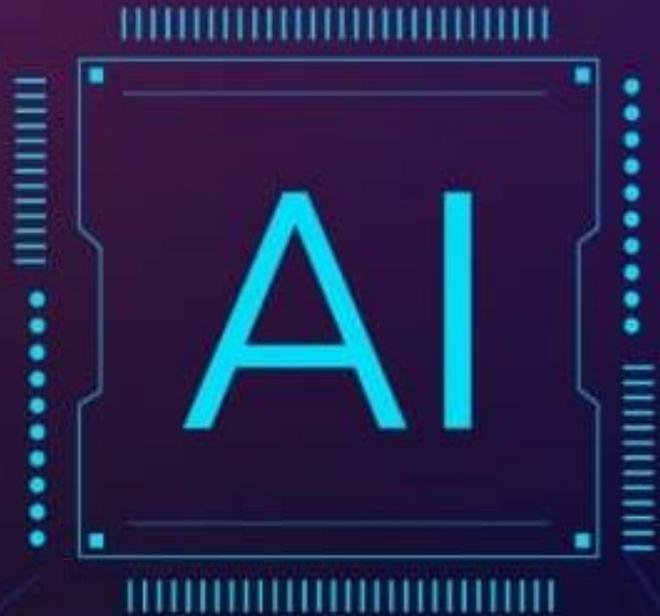
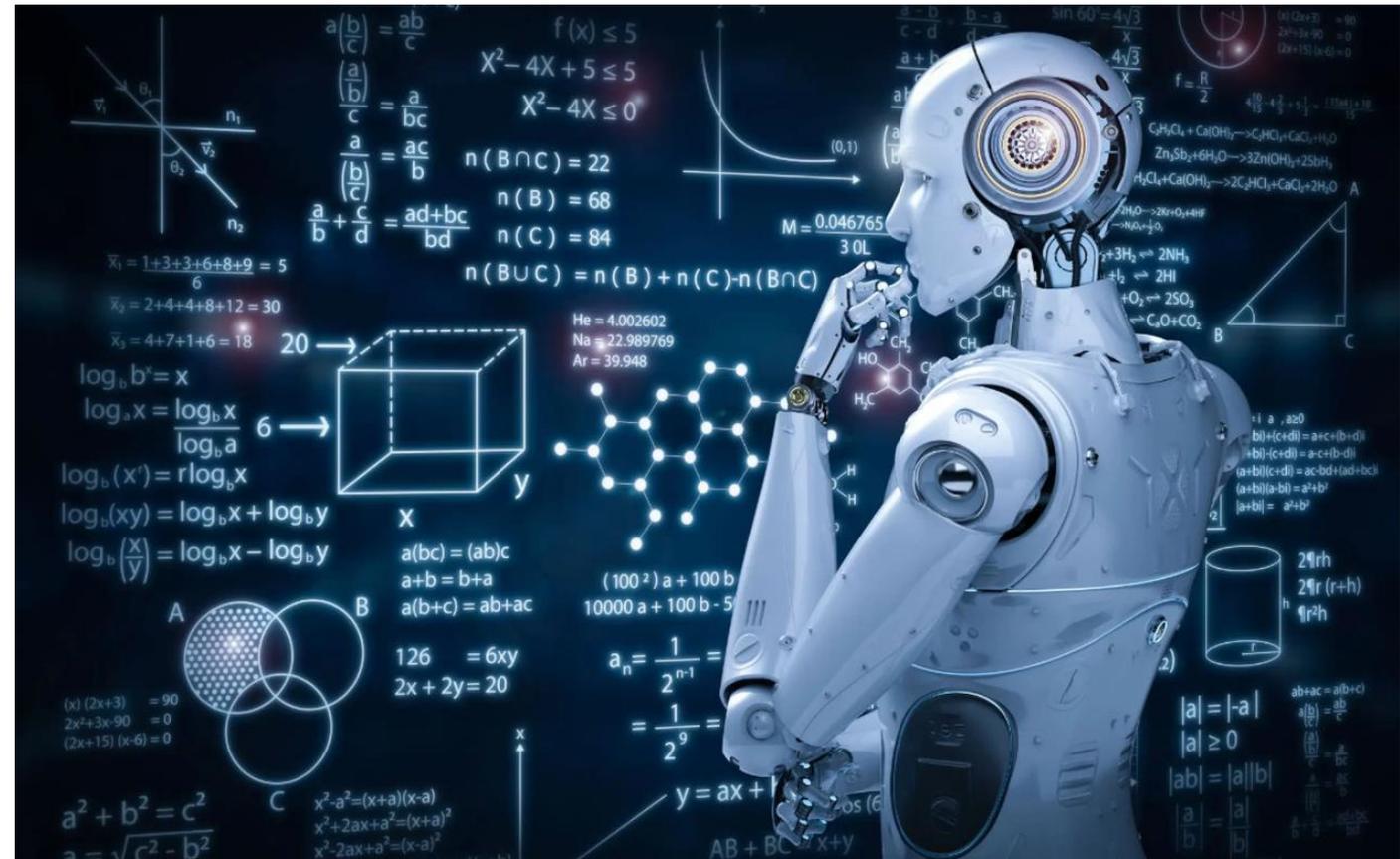


Inserción de la IA en Centros de Datos



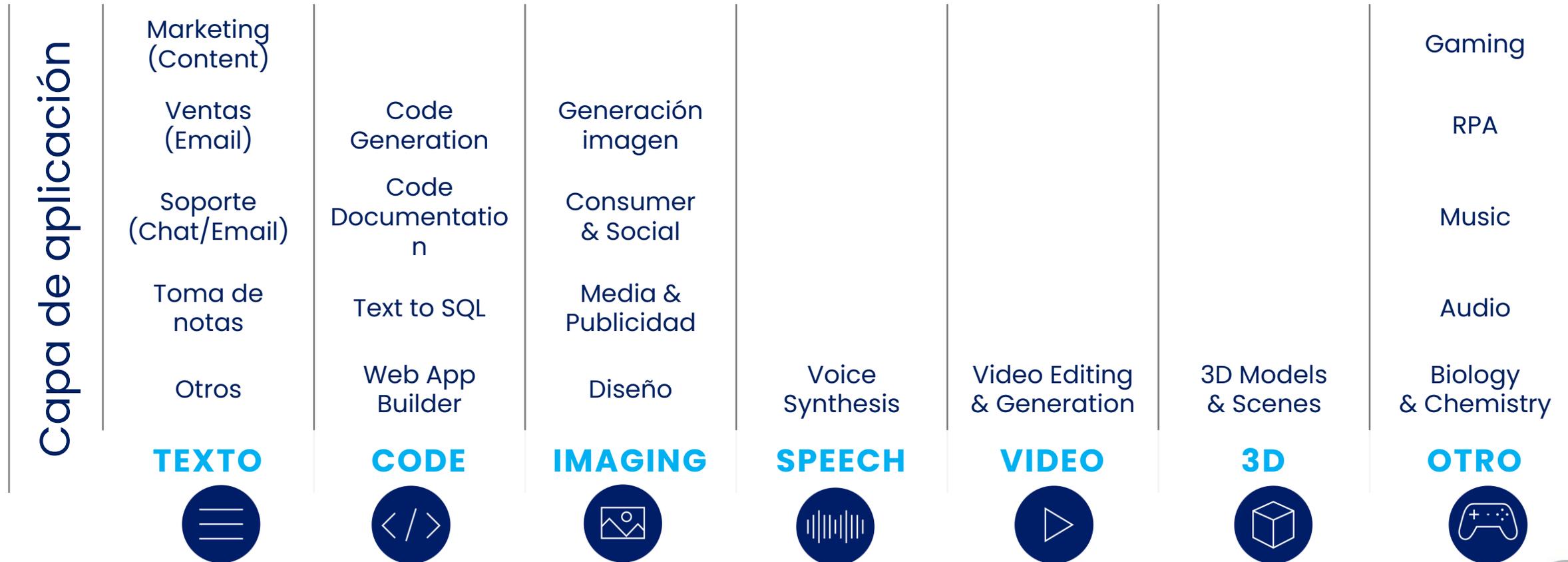
Antonio Cartagena

Redes de Inteligencia Artificial y Machine Learning



Uso de la Inteligencia Artificial en Centro de Datos

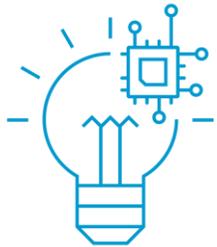
Propósito #1: Manejo de ingresos



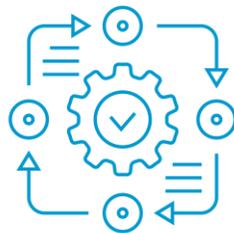
Uso de la Inteligencia Artificial en Centro de Datos

Propósito #2: Mejorar la eficiencia de las operaciones

Energía
Eficacia
& Sostenibilidad



Activo
Rendimiento
Administración



Gestión de la
capacidad
& Planificación



Cliente
Relación
Administración



Seguridad

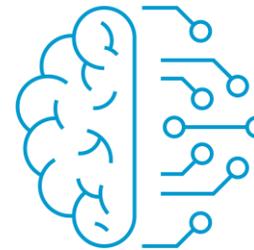
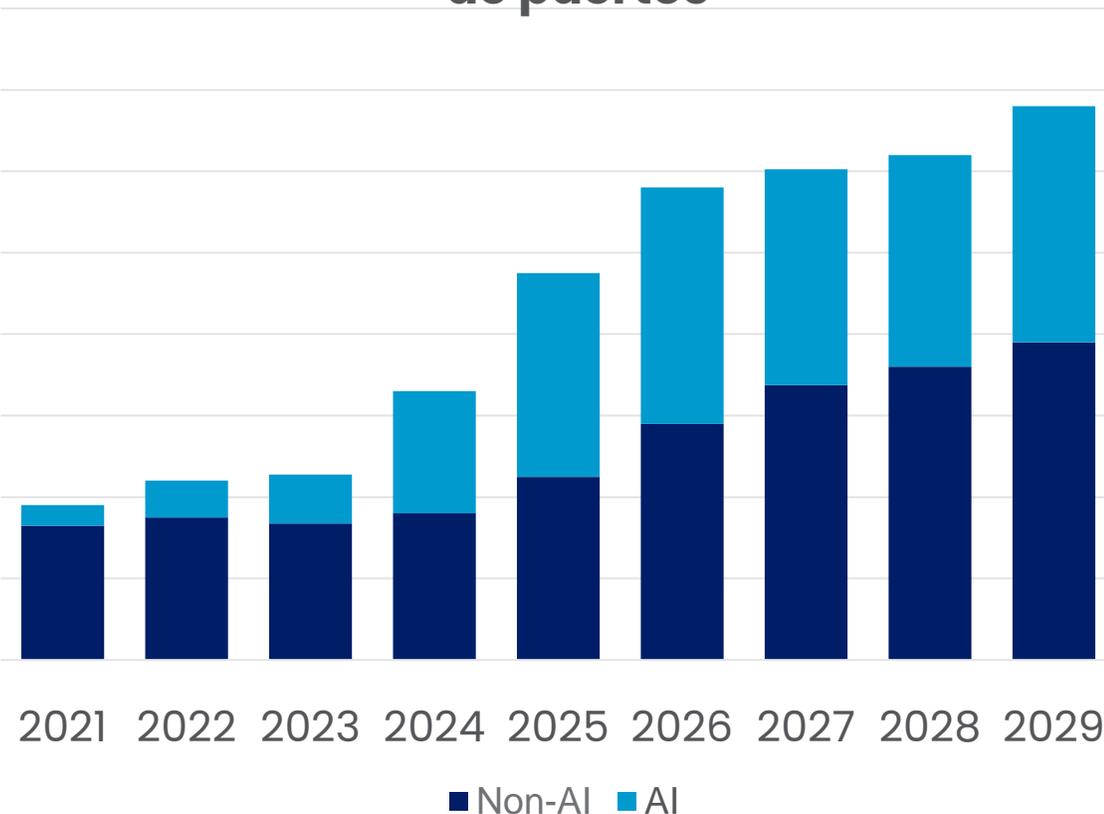


Mejoras
Productividad



Aplicaciones de la IA / ML en el Centro de Datos

Importante impulsor del crecimiento de puertos



AI y Machine Learning crecerá en los Centros de Datos

	% de las ventas de transcievers		CAGR
	2020	2028	2020-2028
AI Clusters	15%	41%	29%
Rest of Cloud Network	85%	59%	9%

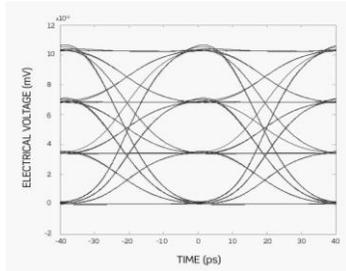
Source: LightCounting

La escala de las implementaciones puede ser muy grande

AI - Directrices de diseño de red

AI XPU Size	Server I/O 100s of XPUs	Rack Scale 1000s of XPUs	DC Scale 10k+ XPUs
AI Network Options	CXL NVLink PCIe	AI Leaf Ethernet or HPS IB	AI Spine IP + Ethernet
	Pequeño AI Apps	Moderado AI Apps	Grande AI Apps

Tres caminos hacia velocidades de datos más altas



1 Velocidad del carril
100G, 200G, 400G, ...

1

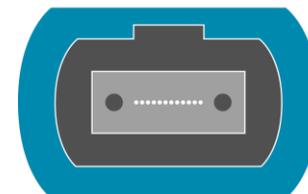
Número de longitudes de onda
1, 4, 8, ...

2



3

Recuento de fibras
2, 8, 16, ...



Etapas de la Inteligencia Artificial

Adiestramiento

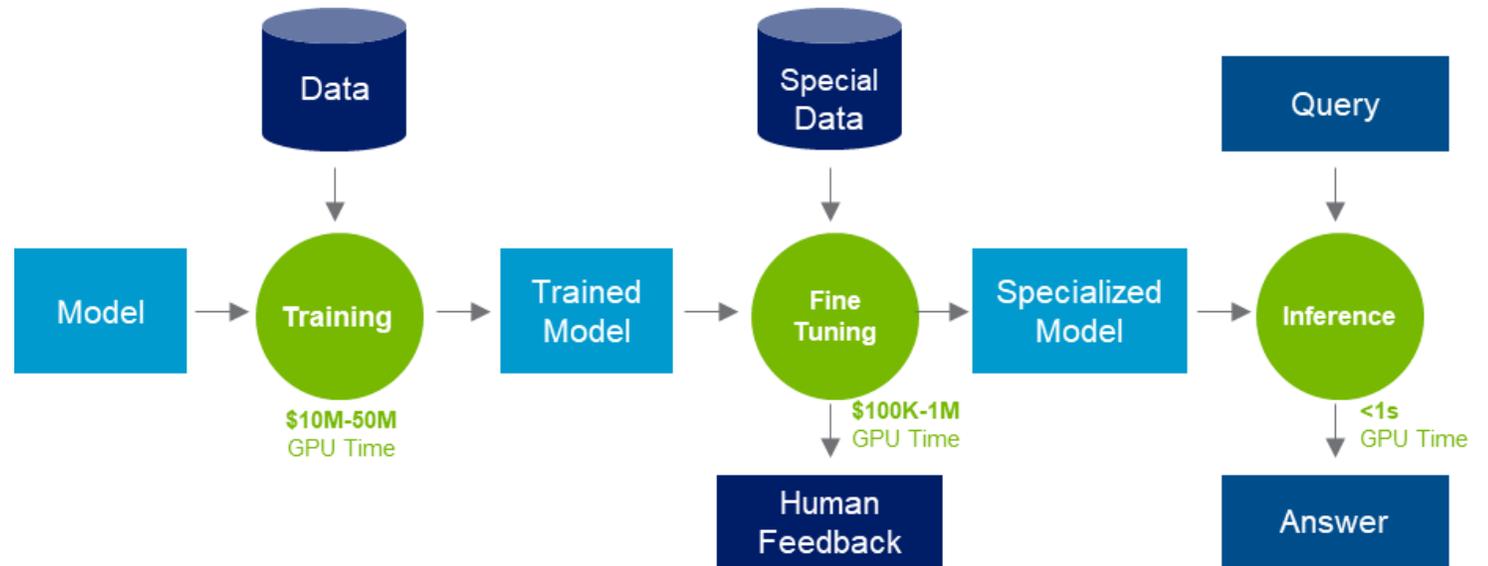
- Intensivo en recursos, mucha potencia durante un periodo prolongado

Ajuste fino

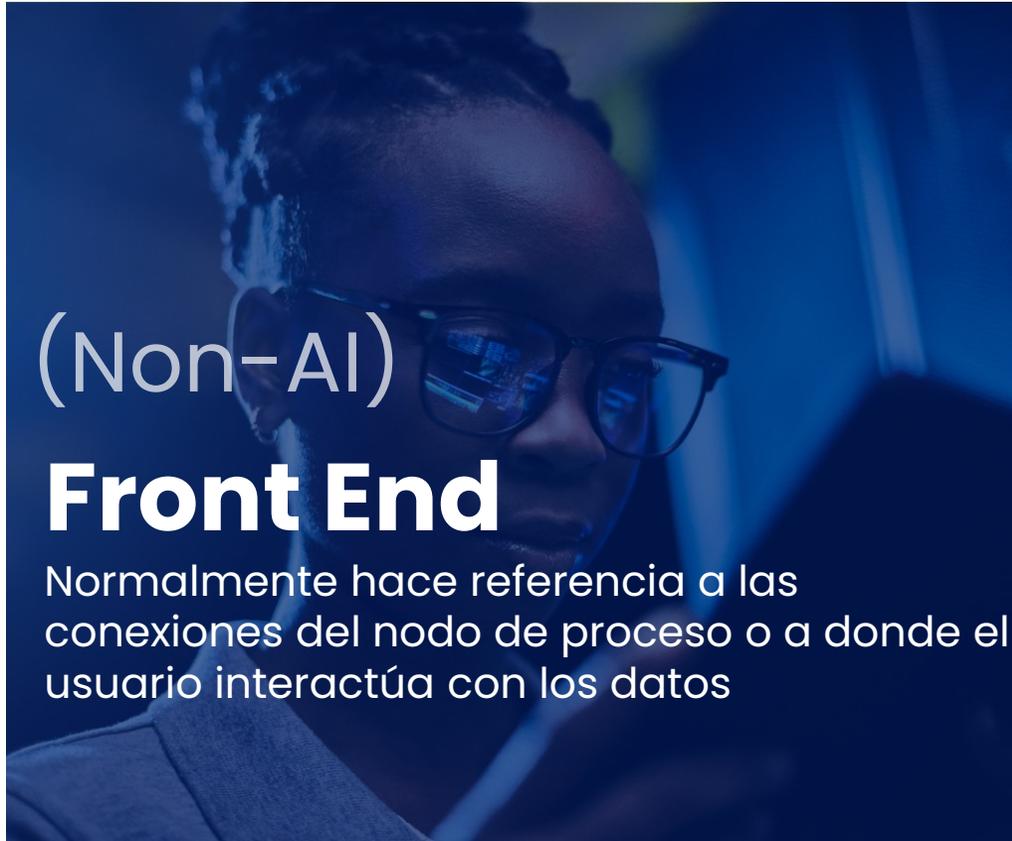
- Utiliza la retroalimentación humana para perfeccionar los datos de aplicaciones específicas

Inferencia

- Uso incremental de energía, respondiendo a consultas individuales



Diferentes partes de la red



(Non-AI)
Front End
Normalmente hace referencia a las conexiones del nodo de proceso o a donde el usuario interactúa con los datos



(AI)
Back End
Por lo general, se refiere a la función de clúster de IA, o donde la mayoría de los procesamientos de datos

Cuestiones que deben abordarse

Energía



NVIDIA DGX H100

“El estándar de oro para la infraestructura de IA”

Datasheet

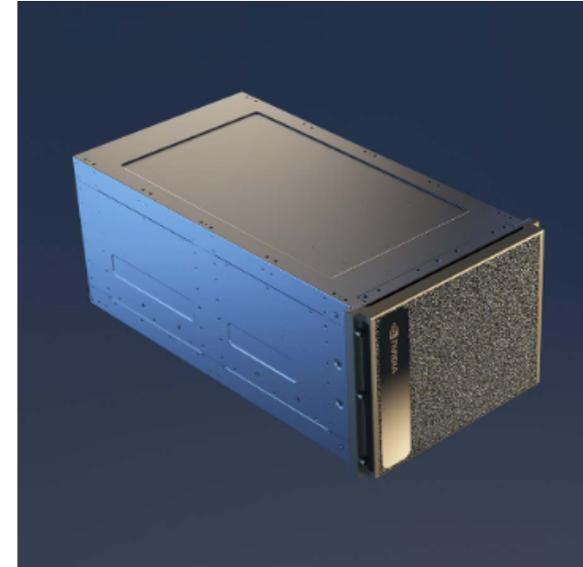


NVIDIA DGX H100

The gold standard for AI infrastructure.

Artificial intelligence has become the go-to approach for solving difficult business challenges. Whether improving customer service, optimizing supply chains, extracting business intelligence, or designing leading-edge products and services with generative AI and other transformer models, AI gives organizations across nearly every industry the mechanism to realize innovation. And as a pioneer in AI infrastructure, NVIDIA DGX™ provides the most powerful and complete AI platform for bringing these essential ideas to fruition.

NVIDIA DGX H100 powers business innovation and optimization. Part of the DGX platform and the latest iteration of NVIDIA's legendary DGX systems, DGX H100 is the AI powerhouse that's the



Specifications

GPU	8x NVIDIA H100 Tensor Core GPUs
GPU memory	640GB total
Performance	32 petaFLOPS FP8
NVIDIA® NVSwitch™	4x
System power usage	10.2kW max
CPU	Dual Intel® Xeon® Platinum 8480C Processors 112 Cores total, 2.00 GHz (Base), 3.80 GHz (Max Boost)

Tecnología en Transceiver

Debido al corto alcance, los transceivers multimodo no siempre necesitan procesamiento de señal digital (DSP)



Product Datasheet

800G QSFP-DD SR8 Optical Transceiver
PN: VD-8CSR8CP-LP

Product Overview

VD-8CSR8CP-LP is a DSP free parallel 800G SR8 based 8-lane QSFP-DD pluggable transceiver that provides independent serial optical data links up to 8x 106 Gbps using PAM4 modulation format over multi-mode fiber.

Features

- Up to 106 Gbps data rate per channel by PAM4 modulation
- Support 800GAUI-8 electrical interface
- Integrated 850nm VCSEL array and PD array without any DSP or CDR
- Single MPO16 connector receptacle optical interface compliant
- Hot-pluggable QSFP-DD form factor
- Low power consumption (4.5W)
- **Single +3.3V power supply**
- RoHS compliant
- IEEE 802.3cm, QSFP-DD HW Rev6.2, and CMIS 4.0 compliant

Applications

- Data Centers and Cloud Networks
- 800G interconnect

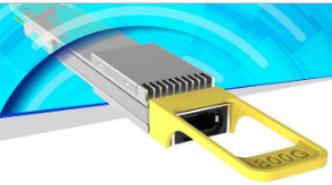


16 W

4.5 W

11.5 W

EOLD-138HG-5H-M
Single-Mode, 800G, 8x100G QSFP-DD
With MPO-16 interface



Product Description

Eoptolink's QSFP-DD 8x100Gbps transceiver module can be used for 800 Gigabit Ethernet connections over 500m of single-mode fiber. The module includes eight parallel channels with a central wavelength of 1310nm, and the operating rate of each channel is 106.25Gbps. These 8-channel PAM4 parallel optical signals can be converted into 8-channel PAM4 electrical output signals; and there are 8 independent electrical input/output channels, which can convert PAM4 electrical input data into 8-channel PAM4 parallel optical signal. The transmitter of the module includes a bi-directional PAM4 re-timer ASIC and 8 EML Lasers. The receiver uses 8 photodiodes and two 4-channel TIA arrays, as well as the PAM4 re-timer.

Features

- Supports 850Gbps
- Single 3.3V Power Supply
- Up to 500m over SMF with KP4 FEC supported at the Host side
- MPO-16 connector
- 8x106.25Gbps (PAM4) electrical interface
- PIN and TIA array on the receiver side
- **Power dissipation < 16W**
- Case temperature range: 0°C to 70°C (commercial)
- Safety Certification: TUV/UL/FDA**
- RoHS Compliant

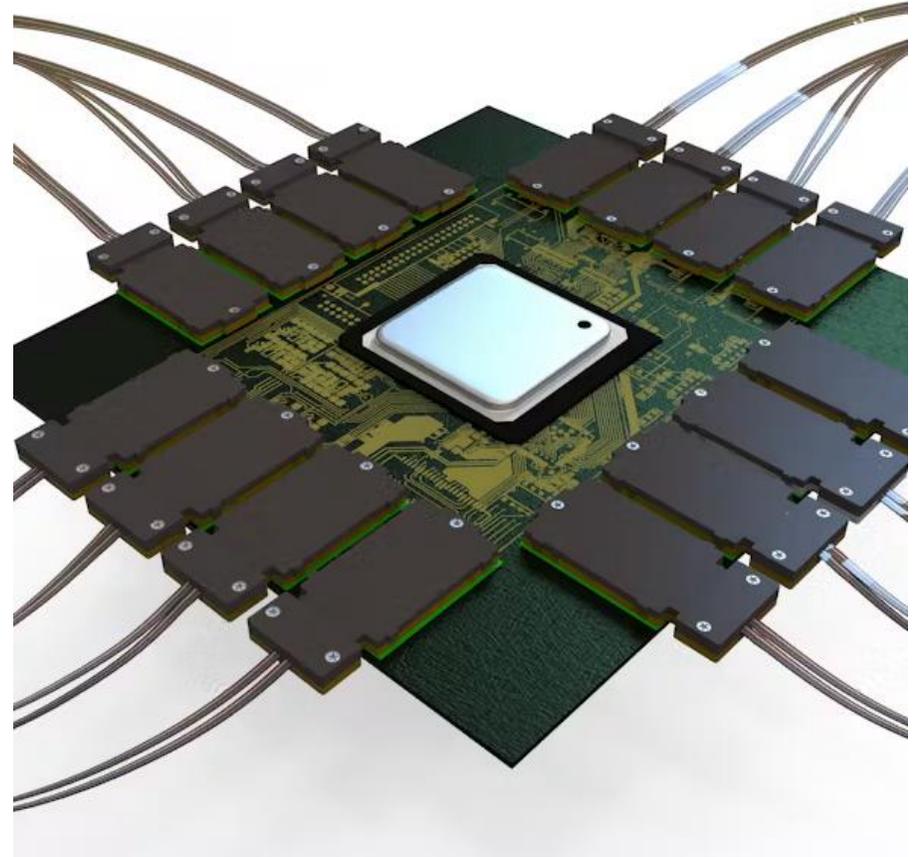
Applications**

- 8x100G Ethernet
- 2x400G Ethernet
- 1x800G Ethernet

Tecnologías futuras en desarrollo

Óptica co-empaquetada

- Promete aumentar la densidad en la cara del interruptor
- Son muchos los obstáculos técnicos que necesitan que se resuelva antes de que esté listo para despliegue masivo





La empresa multinacional de tecnología
NVIDIA utiliza la infraestructura de fibra y cobre de Leviton en sus oficinas y centros de datos en todo el mundo.

20,000

Productos almacenados para NVIDIA en **EE. UU., Reino Unido, Dubái y Hong Kong.**

Geografía

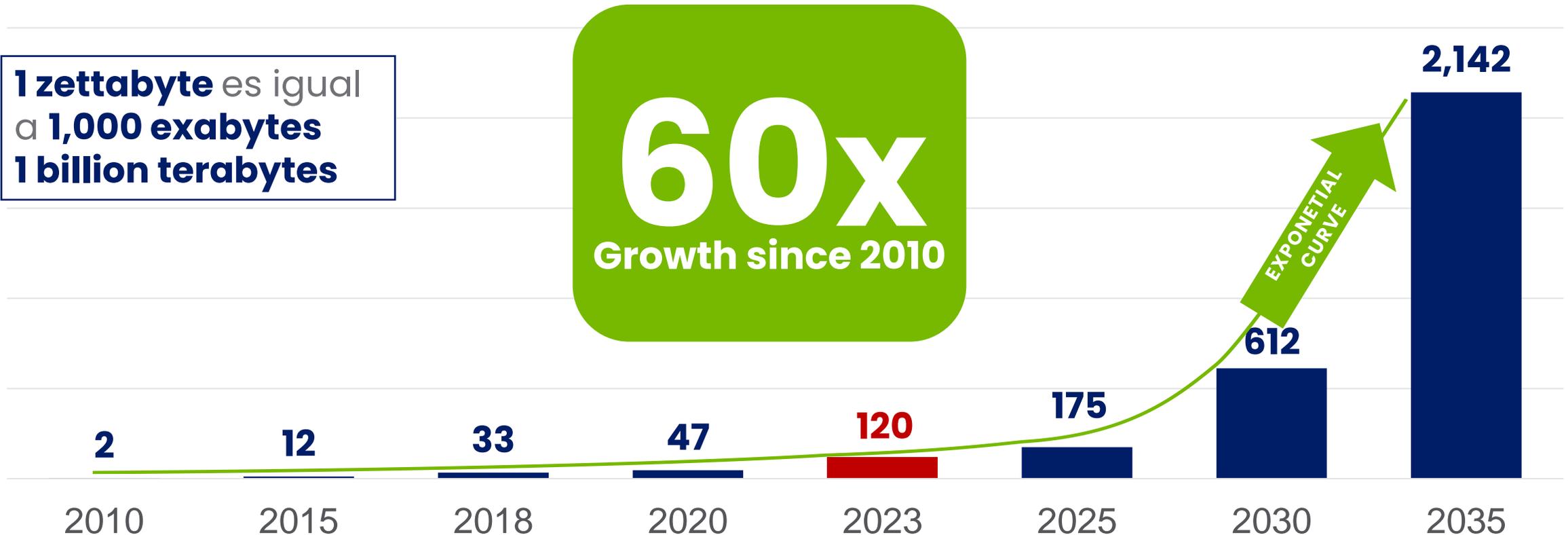


Creación global de datos está a punto de explotar

Cantidad real y prevista de datos creados en todo el mundo en zettabytes

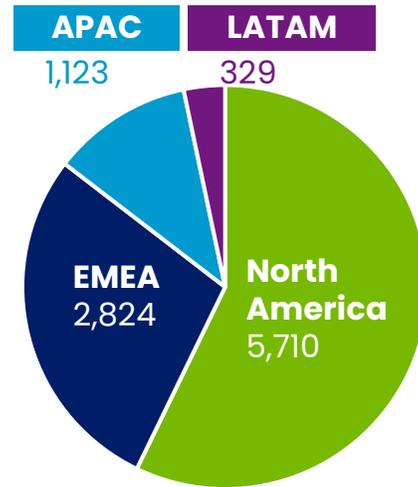
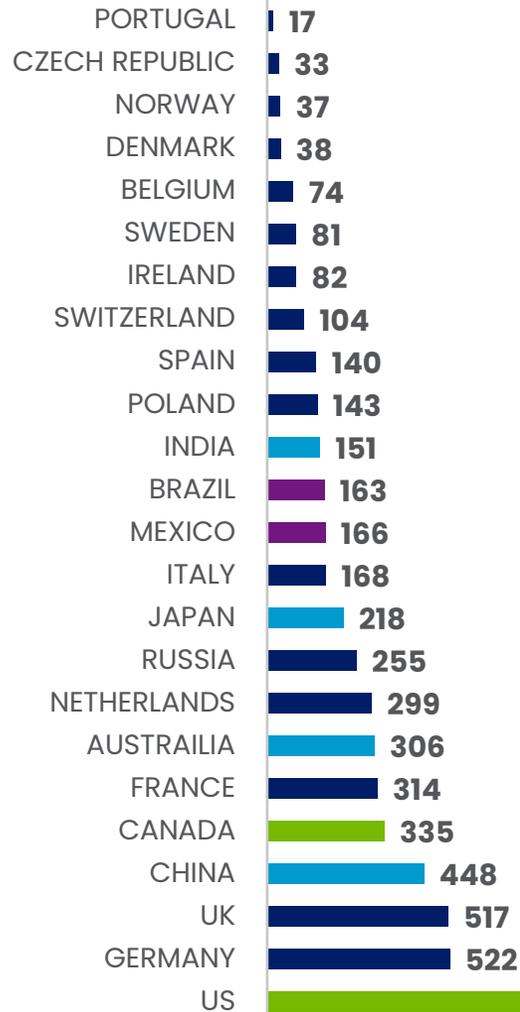
1 zettabyte es igual
a **1,000 exabytes**
1 billion terabytes

60x
Growth since 2010



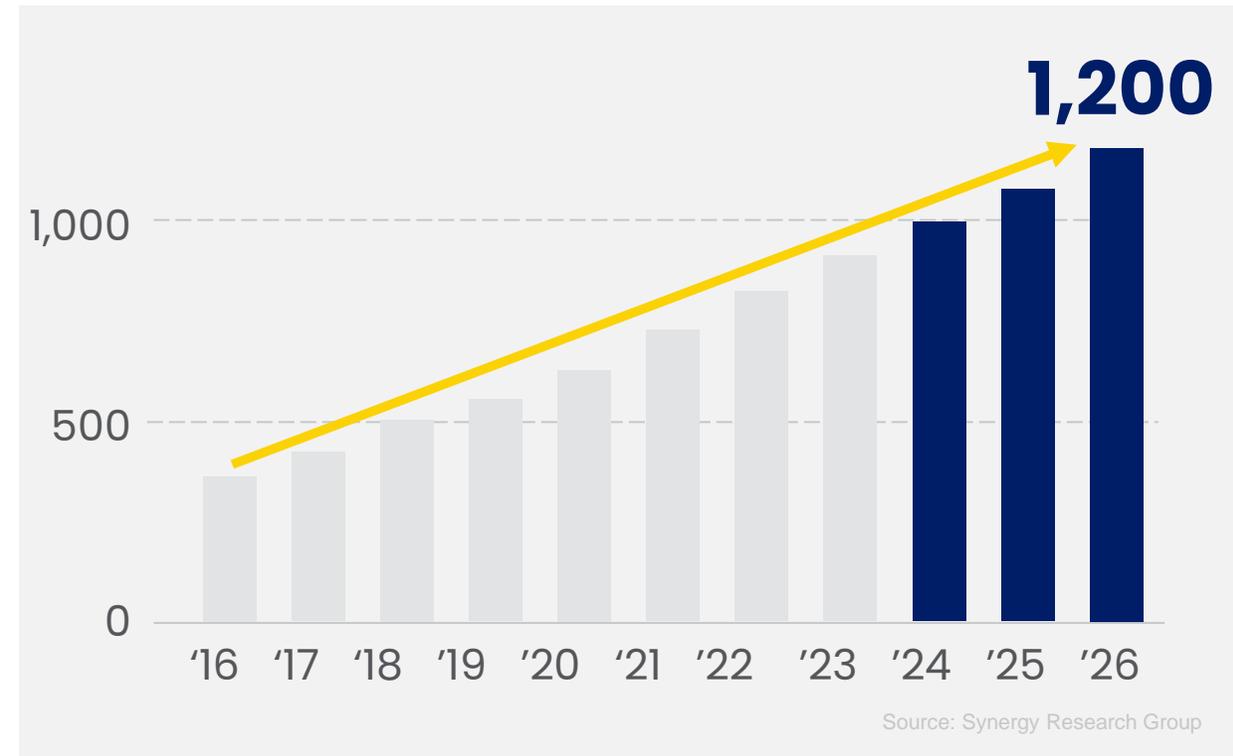
Source: Statista

Centros de datos en todo el mundo



>11K
A partir de
2024
Y creciendo

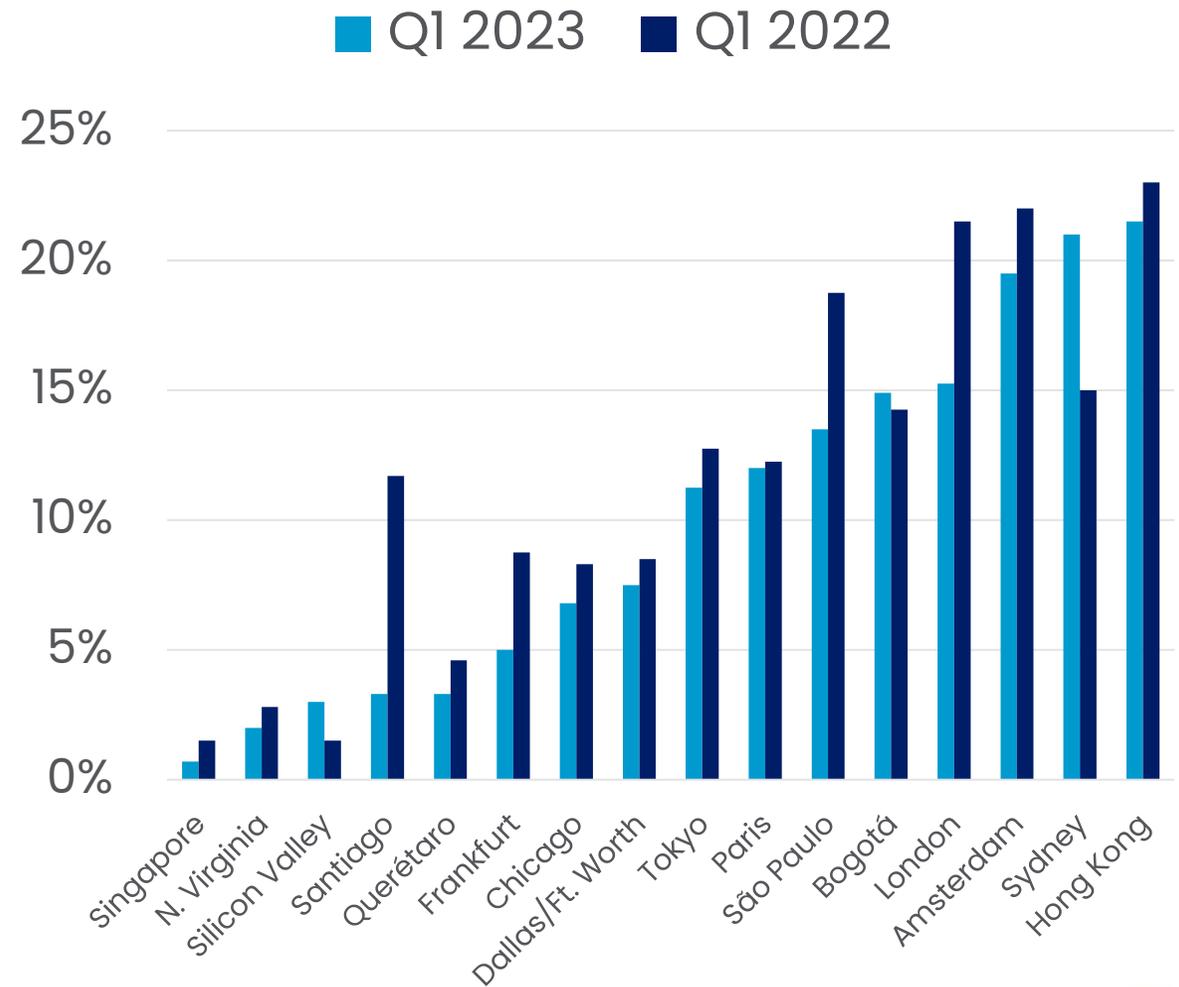
Centros de datos de Hiperescala



Récord de baja desocupación de centros de datos

Tasas de vacantes 2023 VS 2022

- En su nivel más bajo en una década en todos los principales mercados de América del Norte
- La disminución más significativa de Chicago de 8.2% a 6.7%
- Silicon Valley cerca de un mínimo histórico del 2.9%
- La tasa promedio en los mercados FLAP-D cayó un 4.3% hasta el 12.7%
- En 2024 habrá más oferta, aunque se espera que la tasa de vacantes siga siendo baja ya que, la demanda seguirá siendo fuerte



Se necesitan nuevos centros de datos para alojar clústers de IA

Las limitaciones de energía del centro de datos envían la IA a todas partes

La mayoría de los expertos están de acuerdo en que no hay suficiente electricidad sin utilizar para realizar el futuro procesamiento de IA en centros de datos de hiper escala e instalaciones de colocación a medida que se dispara la demanda de GenAI.

Meta construirá
\$800 millones
1er. CD de IA en
Indiana

Gigantesco acaparamiento de tierras de la IA

Los chatbots pueden vivir en la nube, pero funcionan con enormes cajas de hormigón, y están llegando a una ciudad cerca de ti.

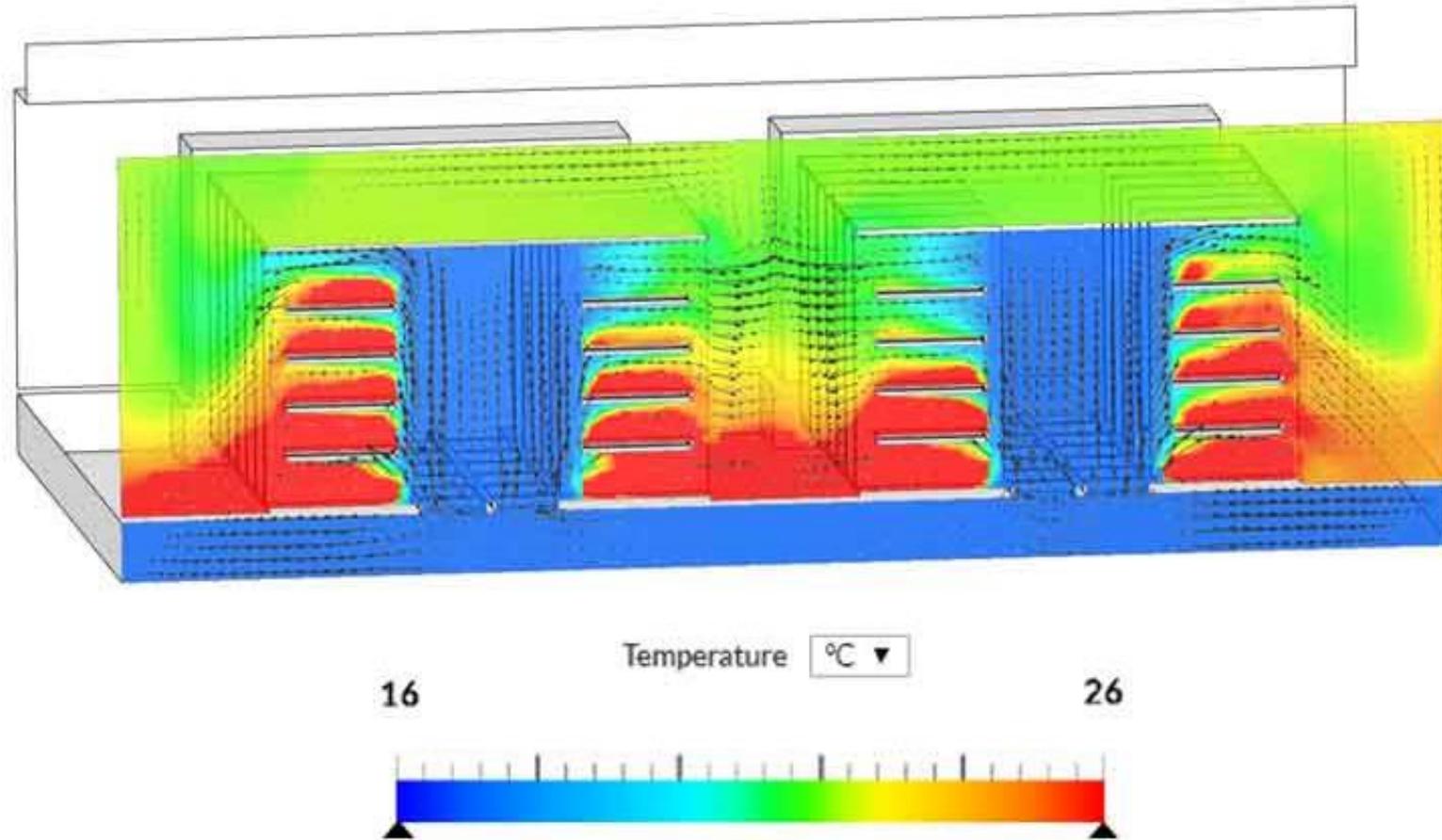
La Inteligencia Artificial requiere **Nueva infraestructura de CD**

La transformación digital de la IA necesitará un salto cuántico en la capacidad de procesamiento.

Enfriamiento



Consumo de energía de la IA



Refrigeración líquida

Asistida por líquido

- Impacto mínimo en el cableado

Enfriamiento directo al chip

- Impacto mínimo en el cableado

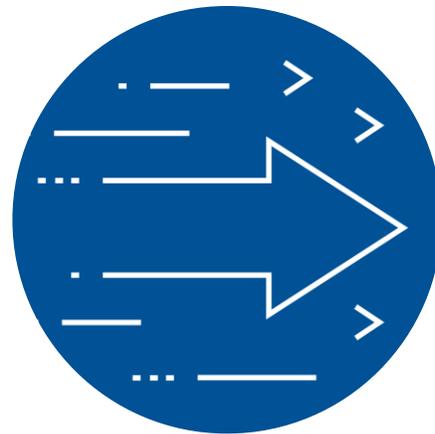
Inmersión en fluidos

- La compatibilidad de los cables es una preocupación



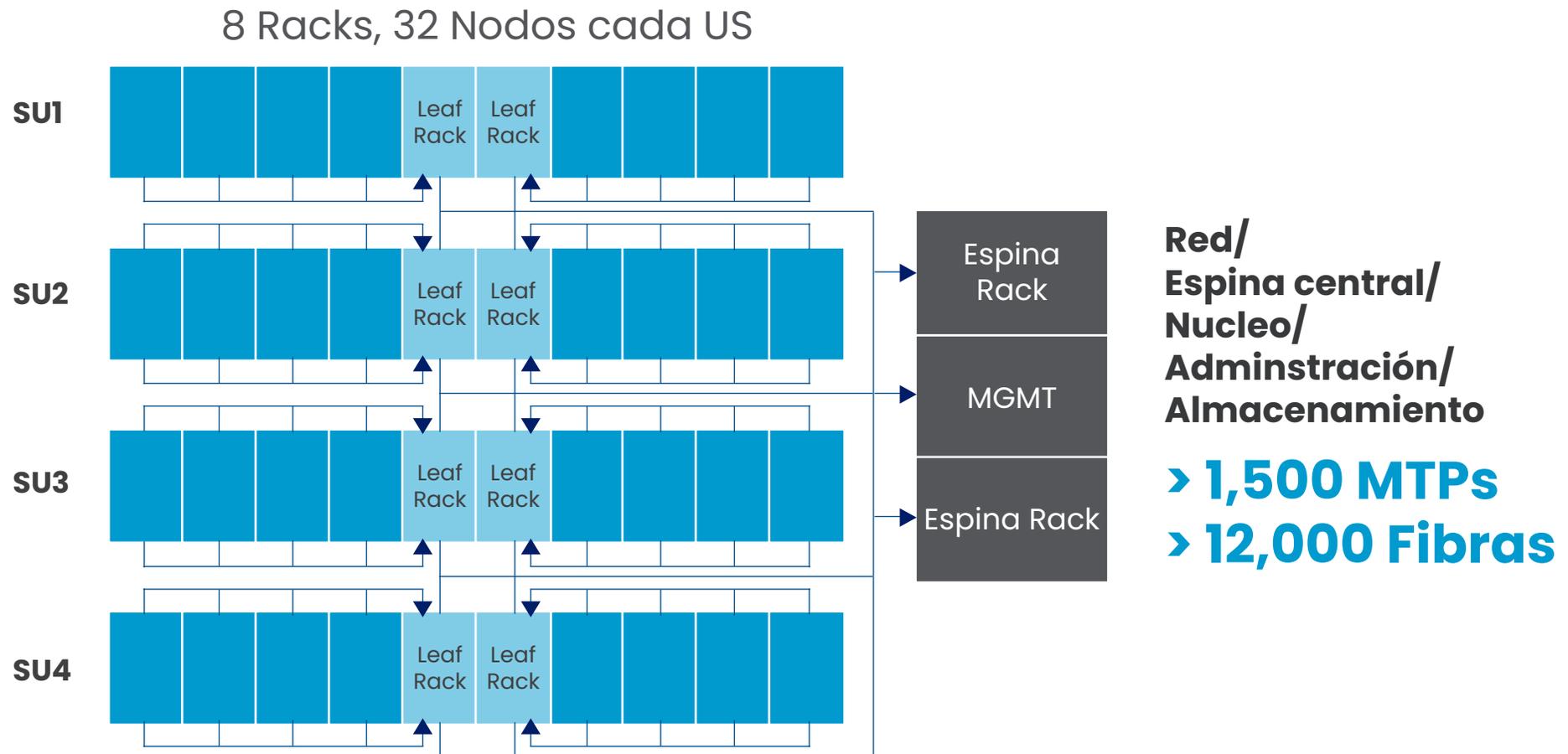
Green Revolution Immersion Cooling

Velocidad de Implementación



Ejemplo de NVIDIA DGX SuperPod

Sistema Completo de 4 Unidades Escalables



IA: Velocidades de datos de red e interfaces con conexión

Las interfaces 400G+ dominan las redes de IA

% of Transceiver Ventas

	2020	2028
AI Clusters	15%	41%
Resto de la Nube	85%	59%

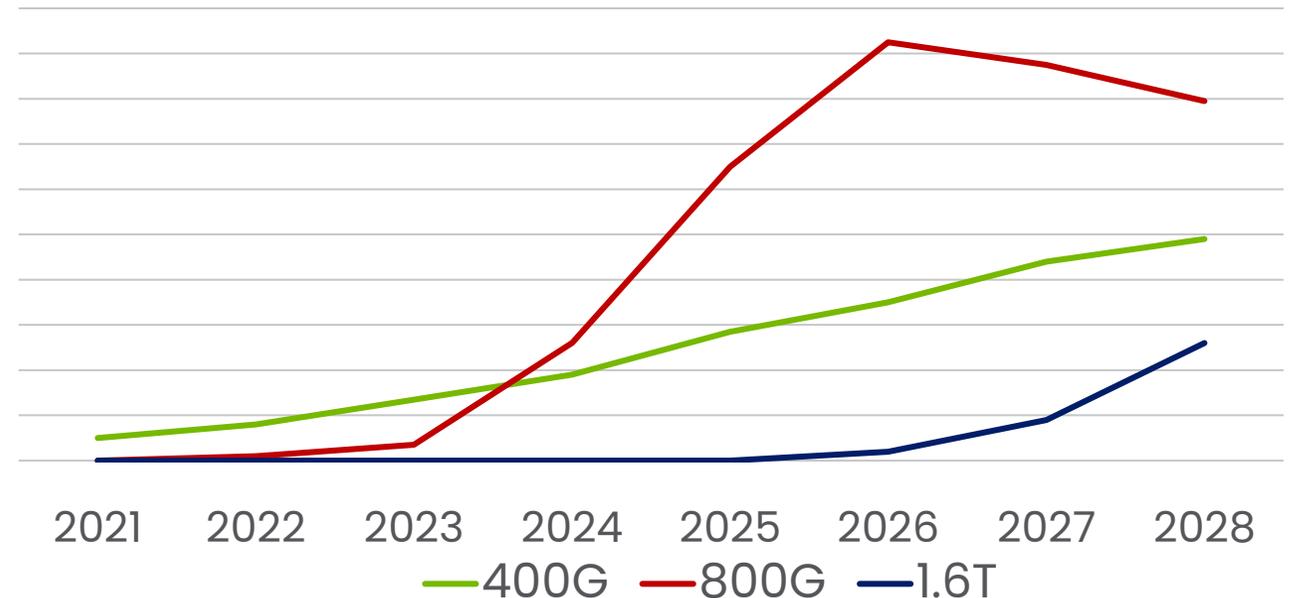
Volumen más alto en transceivers

400G DR4, SR4

800G 2xFR4, DR8, SR8

1.6T FR8, DR8

Volumen de Transceiver ópticos de IA



Las interfaces 400G+ dominan las redes de IA

Publicado 2022

- **IEEE 802.3ck** especifica 100 Gb/s, 200 Gb/s y 400 Gb/s
- Interfaces eléctricas basadas en señalización de 100 Gb/s

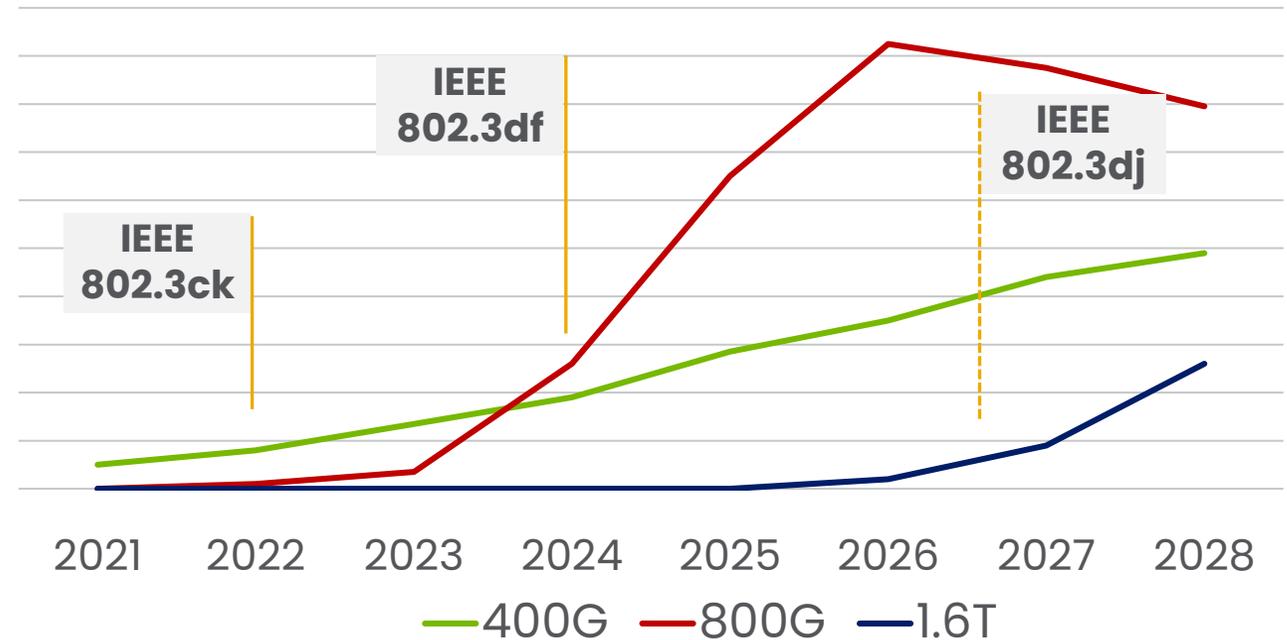
Aprobado recientemente

- **IEEE 802.3ck** especifica 400 y 800 Gb/s con señalización de 100 Gb/s
- Aprobado en febrero de 2024

Actividad actual

- **IEEE P802.3dj** especificará 200-1600 Gb/s con señalización de 200 Gb/s
- Quedan muchos obstáculos técnicos por resolver

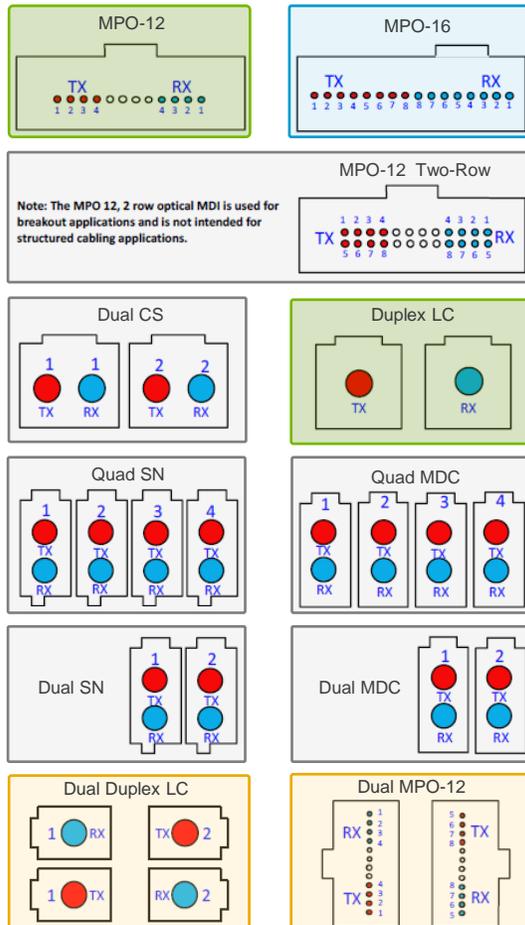
Volumen de Transceiver ópticos de IA



Source: LightCounting

Opciones de conectividad de la unidad MSA

QSFP-DD, -DD800 (Sucesión ahora)



- 400G y conectividad anterior dominada por LC y MPO-12 (FR4, DR4, SR4)
- Las conexiones multimodo se vuelven angulares

- Estos conectores siguen siendo los más populares a 800G, pero desplegado con 100G / carril, vientre con vientre

- MPO-16 se vuelve más popular a 800G, utilizando 802.3df a 100G / carril

- Existen otras opciones basadas en dúplex y MPO-24 con conectores VSFF, pero en un volumen significativamente menor

Arquitectura de la red de IA y requisitos de cableado

Componentes de la red de IA

Cuatro funciones distintas que requieren conectividad

NVIDIA DGX H100 ejemplo:

Compute Fabric

Conexiones de mayor ancho de banda habilitación de la comunicación entre GPU a través de los nodos para actuar como una gran supercomputadora para el entrenamiento y el aprendizaje intensivos de la IA

Storage Fabric

Proporcione **acceso inmediato a los datos compartidos** entre nodos en apoyo de la función de formación y aprendizaje

In-Band management Network

Enlaces de alta velocidad para conectar todos los servicios que administran el clúster

Out of Band Management Network

Conexiones de cobre de baja velocidad para otras funciones básicas de gestión; se conecta a servidores, conmutadores, PDU, etc.



InfiniBand Switch



InfiniBand Switch



Ethernet Switch

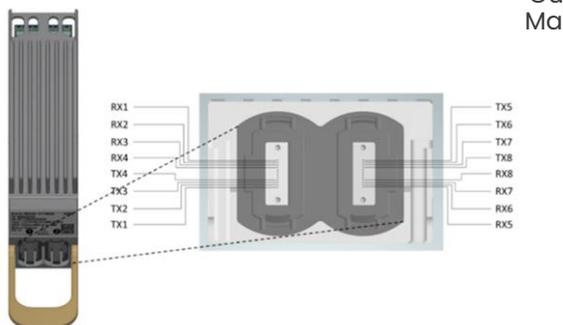
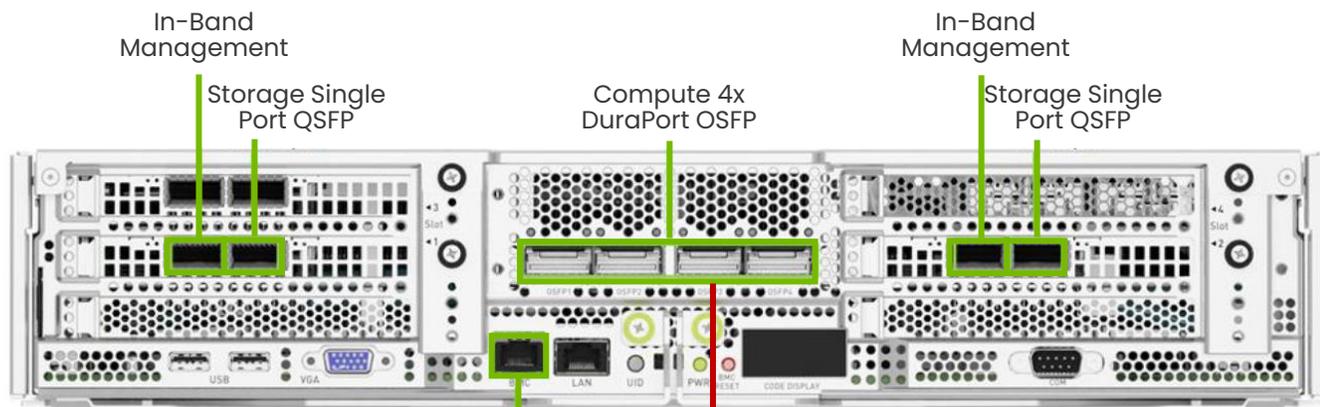


Ethernet Switch

Vista anterior de equipo de red NVIDIA DGX H100

Requisitos de puerto y fibra

DGX H100 Puertos de red



Puertos gemelos transceivers (2x400Gbps) para Conexiones computadas



Nueva interfaz: 8F MM APC

Por servidor requisitos

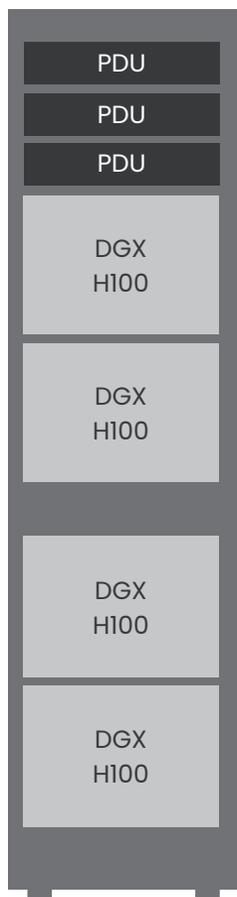
Función	MPO 8 Puertos	Fibras
Procesamiento	8	64
Almacenamiento	1 or 2	UP TO 16
In band	1 or 2	UP TO 16
Total Fiber	10 to 12	UP TO 96

Function	Copper Ports
OOB Cat 6	1

Requisitos de cableado y rack de servidores AI Ileno

NVIDIA DGX H100 Gabinete de servidores

4 Sevidores PorRack



Por servidor Requisitos

Función	MPO 8 Ports	Fibras
Compute	8	64
Storage	1 or 2	UP TO 16
In-Band	1 or 2	UP TO 16
Total Fibras	10 to 12	UP TO 96



Por Rack Requisitos

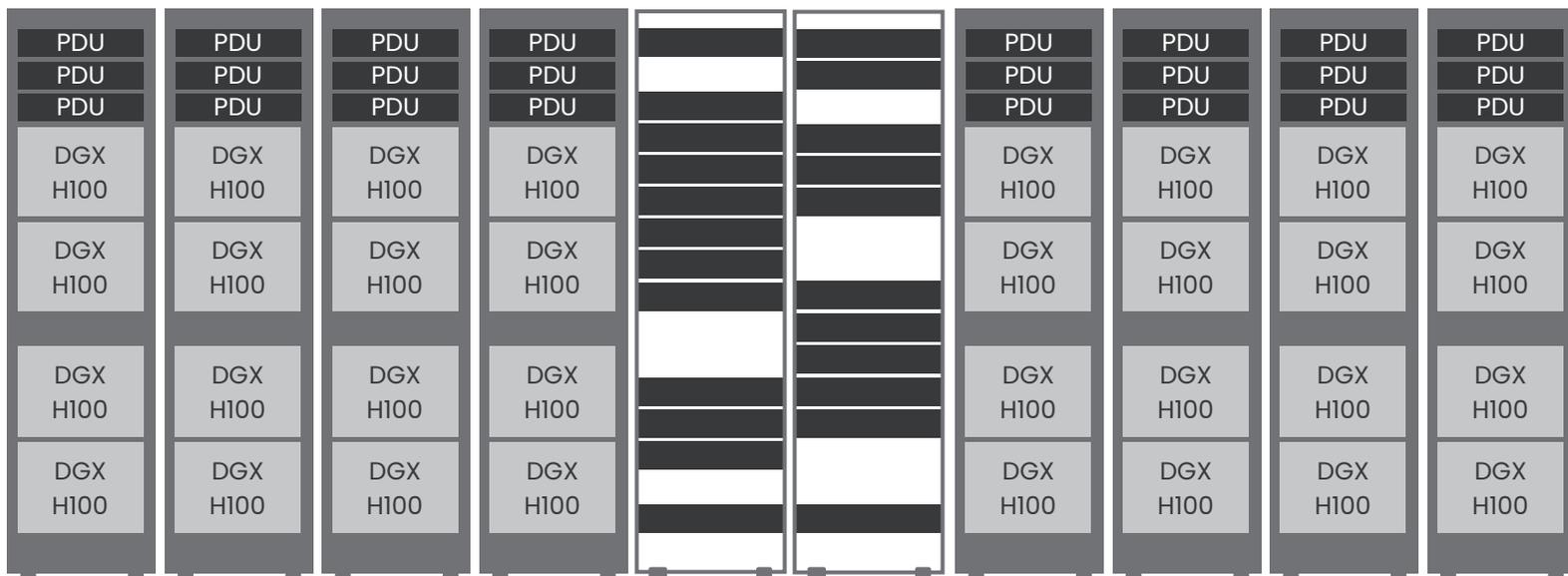
Función	MPO 8 Ports	Fibras
Compute	32	256
Storage	4 or 8	UP TO 64
In-Band	4 or 8	UP TO 64
Total Fibras	40 to 48	UP TO 384

Function	Copper Ports
OOB Cat 6	1

Function	Copper Ports
OOB Cat 6	4

Diseño de "unidad escalable" (SU) de la IA de NVIDIA

Hasta 3,072 Fibras para Switching **hasta 384** MPO to MPO Array Cords, Short Reach, Multimodo



Agregados a Racks Switching

Por rack Requisitos

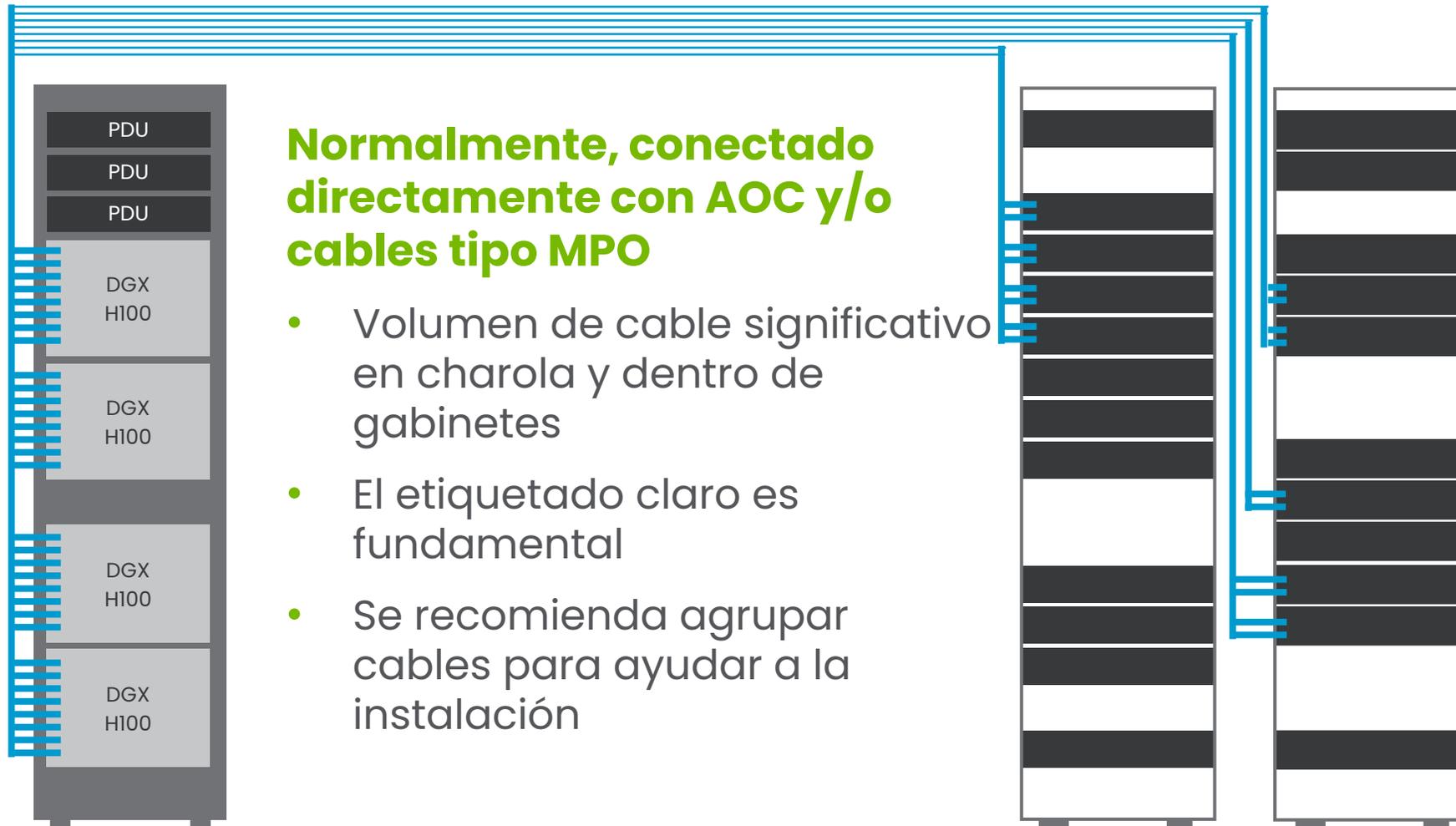
Function	MPO 8 Ports	Fibers
Compute	32	256
Storage	4 or 8	UP TO 64
In-Band	4 or 8	UP TO 64
Total Fiber	40 to 48	UP TO 384
Function	Copper Ports	
OOB Cat 6	4	



Por SU Requisitos

Total Fibras MPO 8	HASTA 384
Total Fibras	HASTA 3,072
Total de cobre nodos	32

Cableado de la red IA

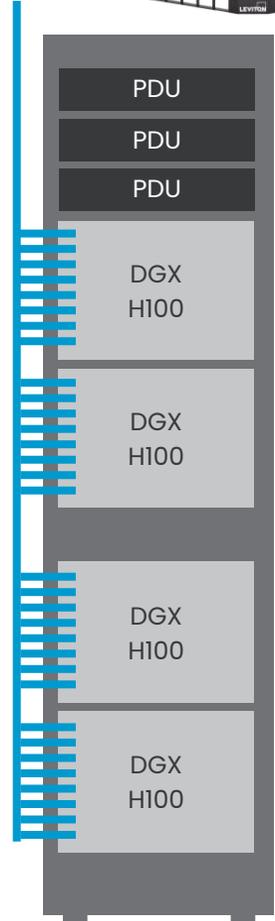


Switching Racks

Cableado de la red de IA – Cableado Estructurado

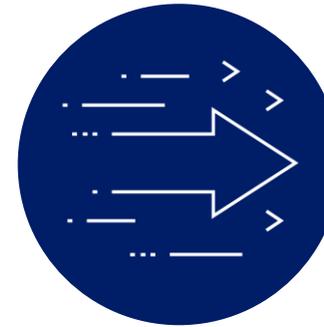


Multifiber Trunks

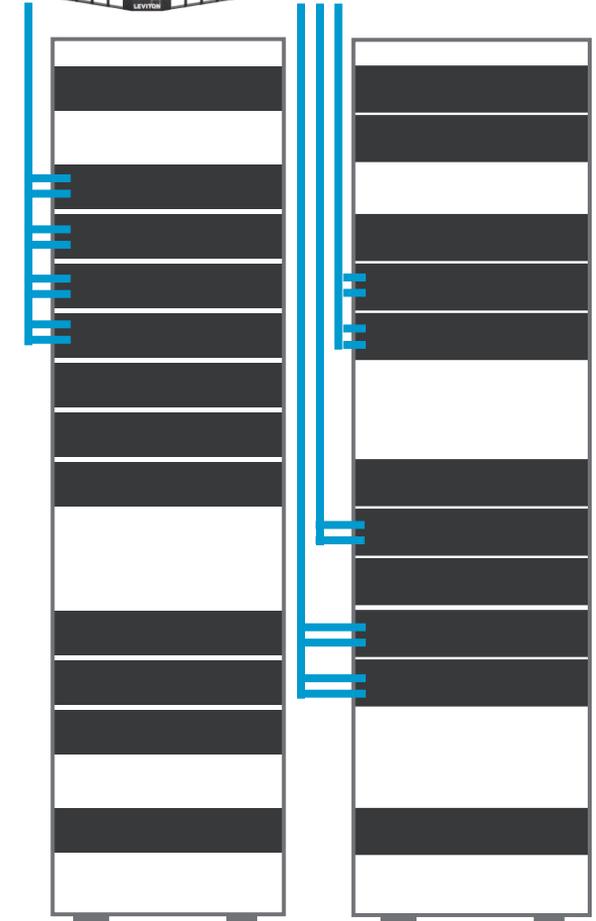


Las soluciones de cableado estructurado tienen muchos beneficios

- Troncales preterminadas multifibra
Reducir la congestión de cables en las charolas superiores
- Más de un 85 % menos de cables que gestionar
- Flexibilidad significativa en la codificación por colores de cables, conectores y puertos
- Soluciones para soportar cobre y fibra en el mismo panel para ahorrar espacio
- Los troncales se pueden instalar antes de que el equipo activo esté en su lugar

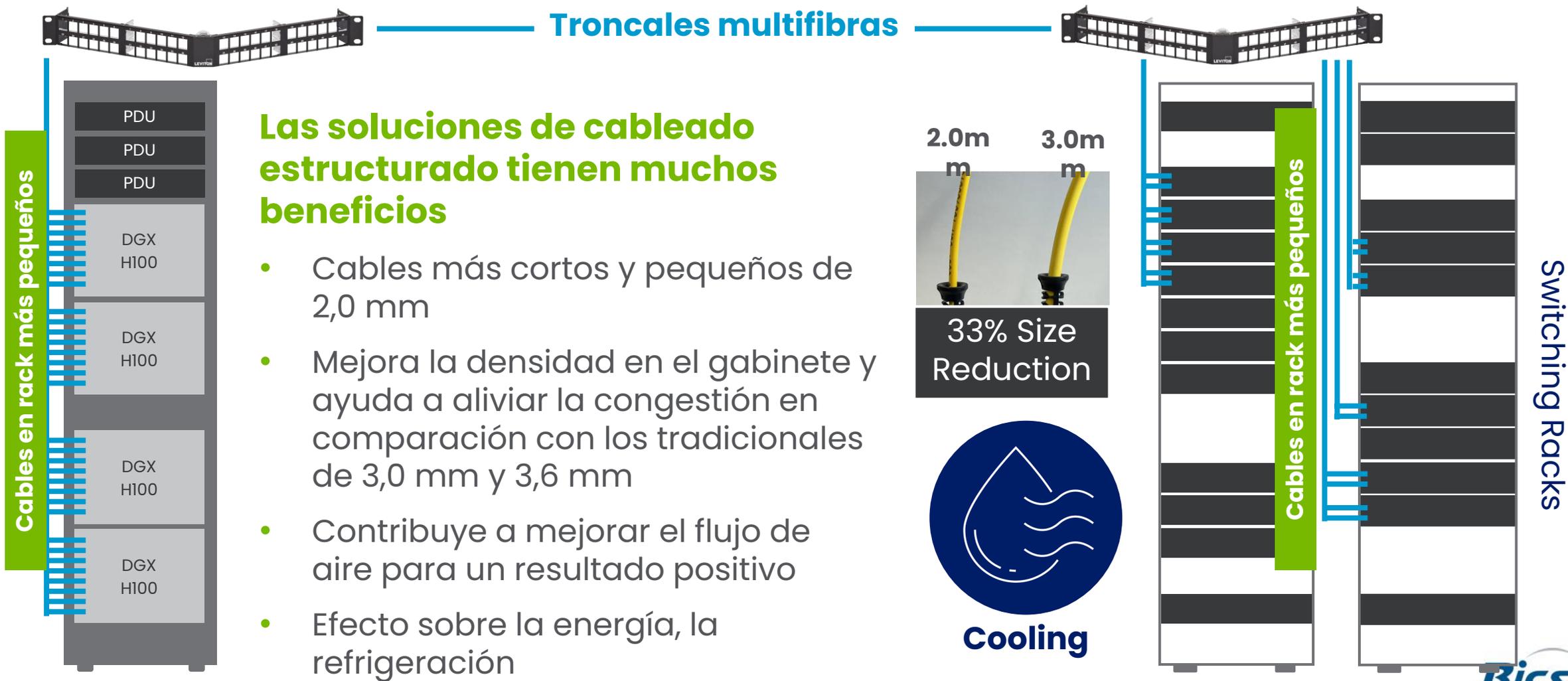


Despliegue
Velocidad



Switching Racks

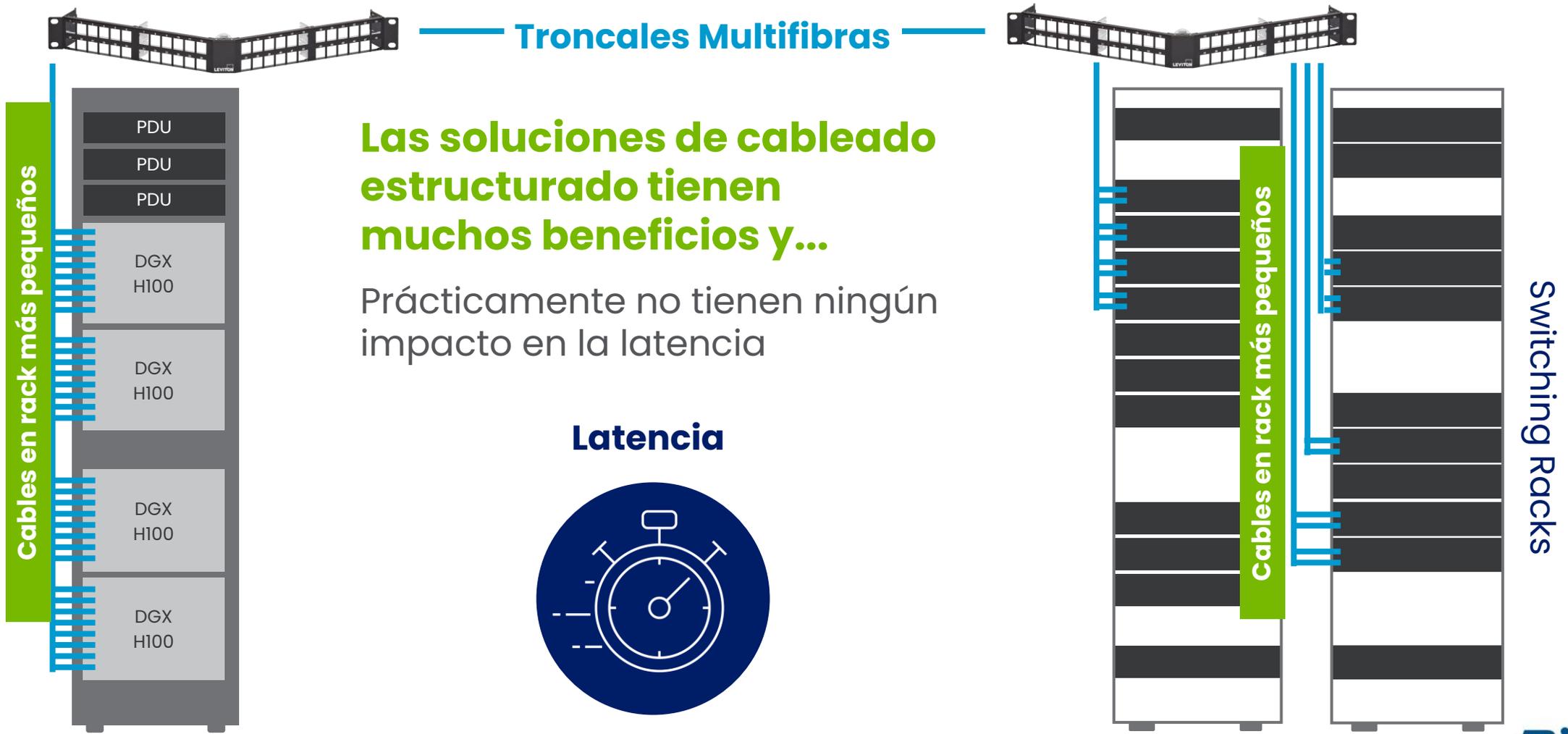
Cableado de la red de IA – Cableado Estructurado



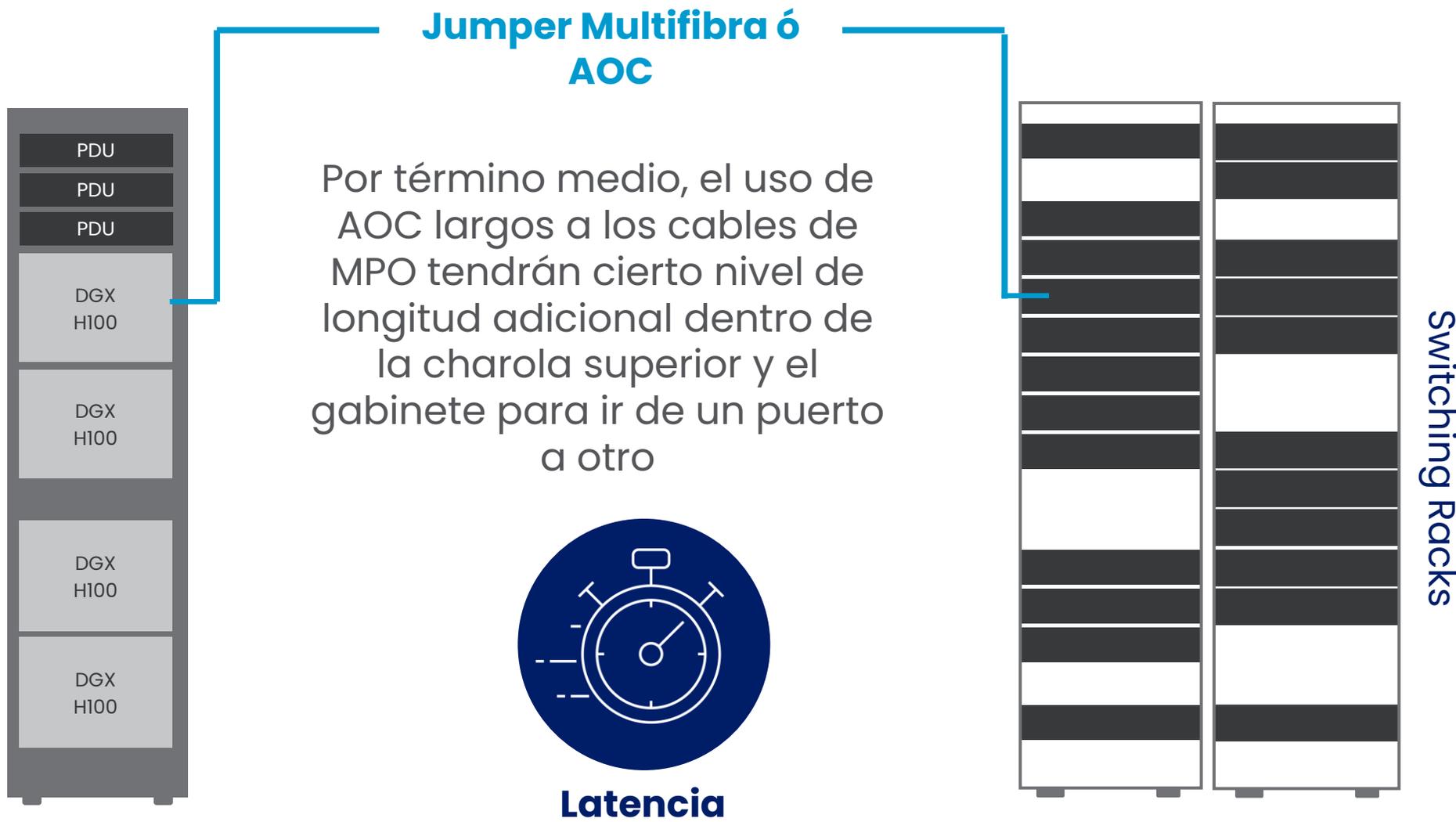
Las soluciones de cableado estructurado tienen muchos beneficios

- Cables más cortos y pequeños de 2,0 mm
- Mejora la densidad en el gabinete y ayuda a aliviar la congestión en comparación con los tradicionales de 3,0 mm y 3,6 mm
- Contribuye a mejorar el flujo de aire para un resultado positivo
- Efecto sobre la energía, la refrigeración

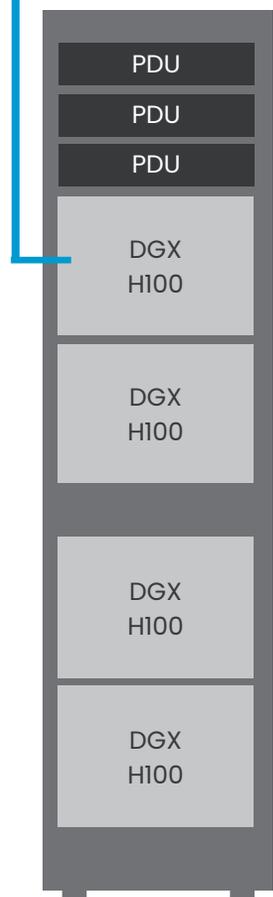
Cableado de la red de IA – Cableado Estructurado



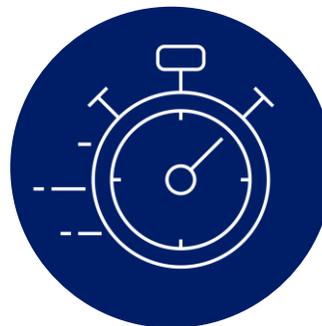
Cableado de la red de IA – Cableado Estructurado



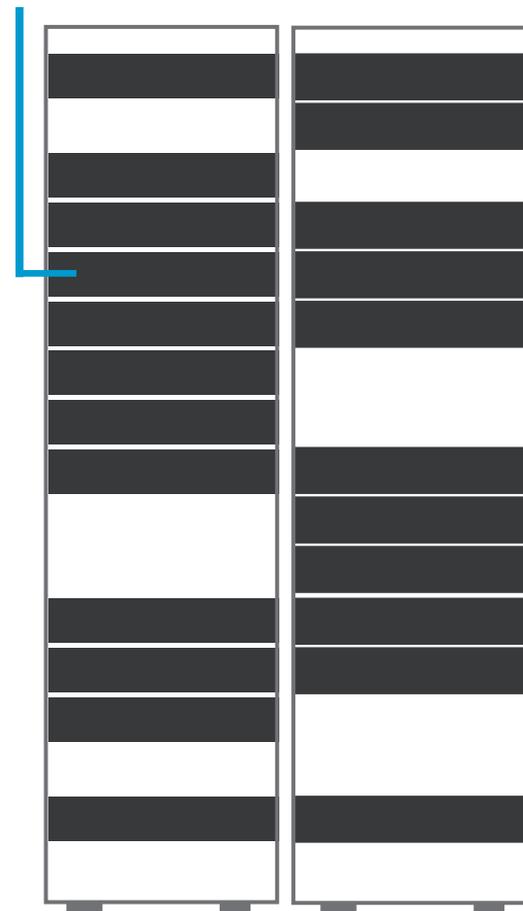
Cableado de la red de IA – Cableado Estructurado



- Los paneles de conexión ya están dentro de la ruta de cable punto a punto deseada
- Las troncales multifibra permiten:
 - Longitudes más precisas en bandejas superiores
 - Jumpers de MPO más precisos dentro del espacio interno
- El uso de paneles fijos o gabinetes con componentes inmóviles minimiza cualquier longitud de cable adicional necesaria

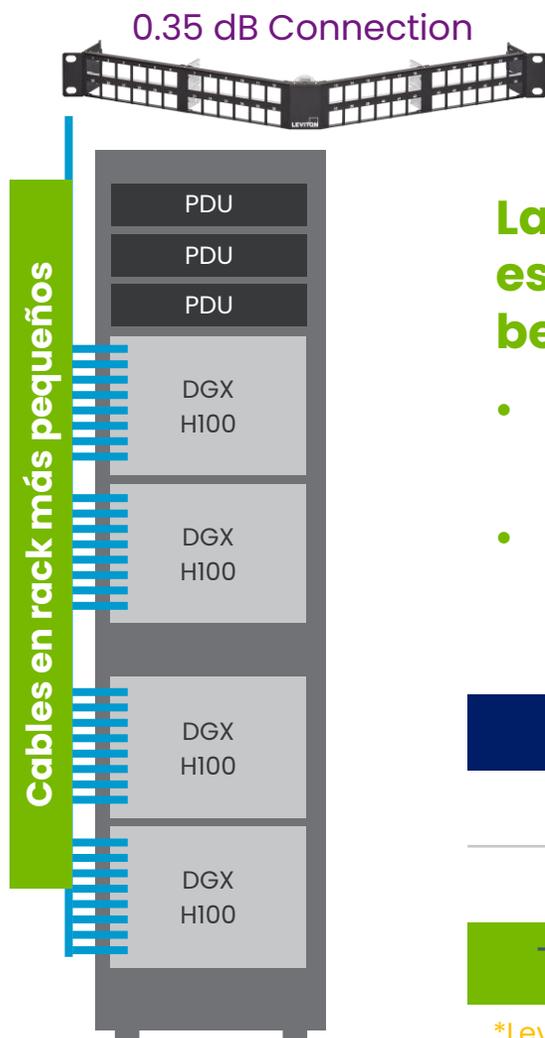


Latencia



Switching Racks

Cableado de la red de IA – Cableado Estructurado



— Troncales Multifibras —

Las soluciones de cableado estructurado tienen muchos beneficios

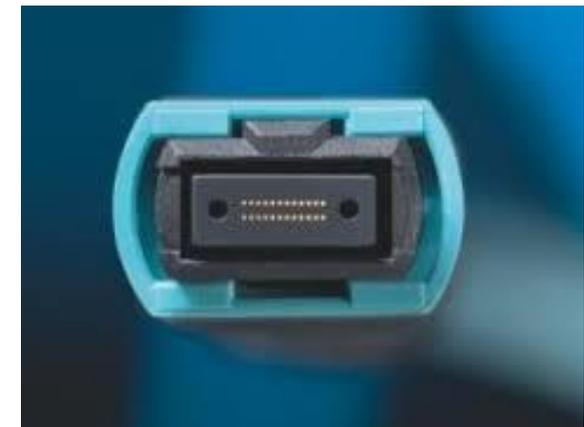
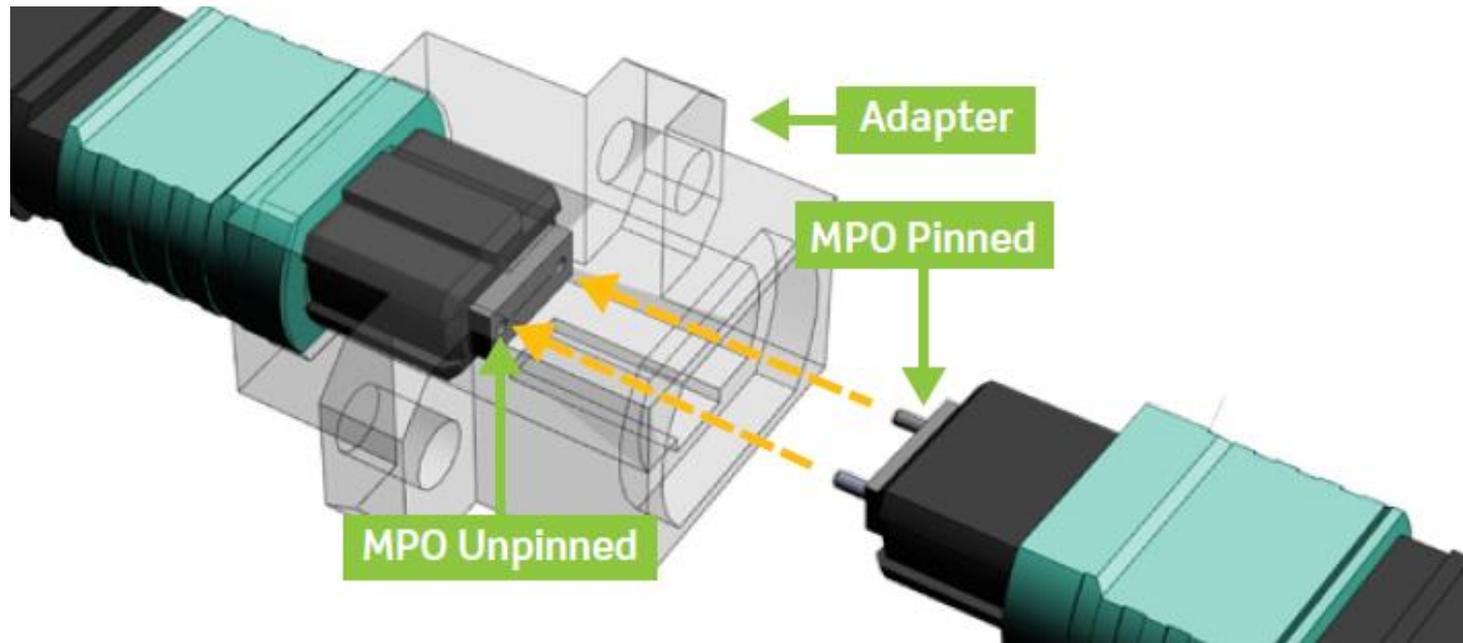
- Prácticamente no tienen ningún impacto en la latencia
- Opere bien dentro del presupuesto de pérdida de enlace

Pérdida	1.8 dB
0.35 dB x2	0.7 dB
10m* Link	0.028 dB
Total IL @ 10m	0.728 dB Max

*Leviton MM OM4 = 2.8 dB/km



Conector MTP y/o MPO



Conector MTP y/o MPO – Códigos de Colores



Cassettes



24-Fiber MTP
(cassette back view)



12-Fiber MTP
(cassette back view)

Multimodo con bota de color en MTP



Red - 24 Fiber



Aqua - 12 Fiber



Gray - 8 Fiber

Monomodo con bota de color en MTP



Red - 24 Fiber



Black - 12 Fiber



Gray - 8 Fiber

Transceivers de 400 Gbps



SFP+
10G



QSFP+
40G



QSFP28
100G



QSFP-DD/OSFP
400G

Arreglo @ 400G: Switch-a-Switch

8-Fiber SM Configuración (10/40/100/200/400G) con DR4

OS2



- Cableado central con base de 24 Fibras
- Conexiones Paralelas (8-fibras) al equipo
- Soporta **400G-DR4/XR4** sobre fibra OS2

Conectividad compatible con aplicaciones 400G+

Conjuntos de productos de Leviton para IA y otras redes de alta velocidad

8F Multimodo y Monomodo



Pulido angulado Array Cord

- La interfaz en ángulo es para transceptores de 400G+
- Conjunto de diámetro pequeño para parcheo de densidad ultra alta y flujo de aire mejorado en el rack
- Podría dividirse en MTP de 1x8F a MTP de 2x4F

8F y 16F MTP placa adaptadora para HDX y E2XHD

- Conjunto completo de placas acopladoras para dar cabida a la tradición Interfaces 8F y 16F emergentes



16F MTP Array Cords and Arneses

- Una fila de 16 fibras, llave de desplazamiento
- Se aplica como interfaz a los transceptores SR8 (MM) y DR8 (SM) 400G, 800G+ 1x 16F
- Agregación y término de fibra hasta APC 2x8F



8F MM/SM APC y 16F MTP Plenum and CPR B2ca Troncales

- Cableado global y estructurado Solución para un mayor número de fibras
- Consolidación en gastos generales bandejas y rejillas



Máxima flexibilidad de implementación

OPT-X™ Sistemas globales de fibra

- Plataformas Alta- y ultra alta densidad
- Arquitecturas flexibles para una variedad de instalaciones
- Engage Low-Loss y Unity Ultra-Low-Loss
- La especificación de cable de fibra multimodo más baja de la industria
- Canales ópticos con rendimiento y alcance más allá de los estándares de la industria



Geografía



Garantía de la aplicación

Experiencia en aplicaciones

- Modelar diseños de red anticipados

Conexiones físicas

- Variar los tipos de fibra y conectores

Transceiver Banco de pruebas

- Comprender la respuesta esperada



Equipo de diseño para Centros de Datos

- Se asocia con los clientes para crear diseños de infraestructura de red eficientes y rentables
-
- Expertos con experiencia en equipos activos en el mundo real y conocimientos de planificación



Resumen

- **Potencia, refrigeración, ubicación geográfica, latencia y velocidad de implementación** son más importantes que nunca en el despliegue de redes de IA
- **Cableado estructurado** puede ayudar a abordar estos requisitos críticos
- Leviton tiene un **conjunto completo de productos** y sistemas globales definidos para satisfacer las necesidades de las redes de IA en todo el mundo



Antonio Cartagena

Gerente de Ventas Network Solutions Centro América y el Caribe

Acartagena@leviton.com

Celular: +1 787 903 0304

¿Preguntas?

The
Industry's
Best
& **Service**
& **Support**

Committed
To Our
Customers
& the
Environment

Outstanding
Return On
Infrastructure
Investment

A **Culture** of
Ingenuity &
Innovation

Quality &
Performance
in every
Solution