SUMMIT

September 9-11, 2025 | Santa Clara Convention Center

2025

DIAMOND PARTNERS



aws cādence Google Microsoft





HIGHLIGHTS

- NVIDIA unveiled Rubin CPX, a new GPU for long-context AI reasoning which enables up to 8 exaflops of NVFP4 compute.
- The age of inference is here. Google revealed its internal inference has exploded: from ~9.7 trillion tokens per month in April 2024 to over 1,400 trillion tokens per month.
- Kove, SK Hynix, Pliops and Cadence showcased significant breakthroughs in the fields of memory, storage and design innovation for large scale Al inferencing systems.
- "Al is kicking our butts and teaching us that we know nothing about infrastructure..." said Meta's VP of Engineering, YJ Song, as he shared their massive datacenter projects including the Prometheus Supercluster.
- Lightmatter's CEO demonstrated unprecedented 800 Gbps bidirectional bandwidth from photonic systems that could spur the next 1000X in AI computing.



ANNOUNCEMENTS AT AI INFRA SUMMIT 2025



EXCLUSIVE NVIDIA ANNOUNCEMENT

Rubin CPX GPU: A New Class of GPU Designed for Massive-Context Inference

Giga-Scale Al Factory Reference Designs

OTHER NEWS

MLPerf Inference v5.1 Benchmark Results announced on-stage

Recogni unveils new brand, and is reborn as Tensordyne

Kove: SDM benchmark results show 5X larger AI inference workloads on Redis and Valley

UALink, RISC-V, CXL Consortium and others showcase products on expo floor



TOP SPEAKERS



VP of Hyperscale and High-Performance Computing NVIDIA



BERTA RODRIGUEZ-HERVAS
Chief Al and Analytics Officer
Pfizer



MARK LOHMEYER

VP & GM, Compute
and Al Infrastructure

Google



BARRY COOKS

VP, Compute Abstractions

AWS



YEE JIUN SONG

VP, Engineering

Meta



LAURA ORTMAN
CEO
Cologix



DAVID GLICK
SVP, Enterprise
Business Services
Walmart

THANK YOU TO OUR 2025 PARTNERS

DIAMOND PARTNERS



PLATINUM PARTNERS









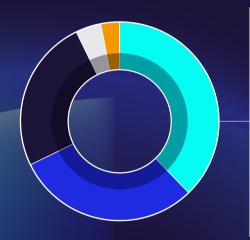




Qualcomm **SIEMENS** SYNOPSYS°

TENSORDYNE

AUDIENCE BREAKDOWN



•	Tech Ecosystem	38%
	Semiconductor, Systems, Systems and Device Manufacturer	
	Hyperscalers, Al Software Infrastructure, Al Application Software	vare,
	Physical Infrastructure	

Enterprise Al Practitioners3

O Hyperscaler......25%

• Investors4%

80%

plan to purchase/upgrade their GPUs in the next 12-18 months. A further 54% expect to do this with memory, and 36% on interconnect 72% annually purchase/upgrade to

new hardware components (e.g memory, chips, interconnect)



ATTENDING COMPANIES INCLUDED...









BERTA RODRIGUEZ-HERVAS, CHIEF AI AND ANALYTICS OFFICER, PFIZER















RAJOSHI GHOSH, CO-FOUNDER & CHIEF ECOSYSTEM OFFICER, PROMPTQL

Top use-cases:

A unified strategic
accurately resolves

Embedded app in S language powered every seller on their

Why it matters:

Unified real-time at









PRESS COVERAGE

GLOBAL MEDIA INCLUDED:















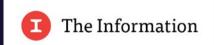








Tech**A**rena®



The New York Times



LOOKING TO 2026

After a watershed year that witnessed a salvo of product launches, including NVIDIA's Rubin CPX, we'll return next year and dissect another year's worth of technology, defining the next 12 months in infrastructure development.

Prestigious keynotes, blockbuster announcements, and the most comprehensive industry representation anywhere in AI infrastructure – that's what makes AI Infra Summit so impactful. 8000 attendees, 400 speakers and 250+ partners will fill the halls of Santa Clara Convention Center in September 2026.

As enterprises enter the era of inference, and Al moves from servers into physical and embodied systems, we will bring you the keynotes, product launches, and technical prowess that you need to stay competitive.

Engage with us now to get your work onto Al infrastructure's ultimate stage.







NVIDIA KEYNOTE ON DAY ONE





September 15-17, 2026
Santa Clara Convention Center

"Exactly the people I want to meet, who I want to talk to and to share our vision, our strategy, our future."

CHIEF INNOVATION OFFICER, ALIGNED DATA CENTERS

POWERING FAST, EFFICIENT & AFFORDABLE AI

- HARDWARE & SYSTEMS
- PHYSICAL AI
- ENTERPRISE AI
- AI DATA CENTER
- 2x DEMO STAGES

8,000 400 250+
REGISTERED SPEAKERS PARTNERS

SPONSORSHIP OPPORTUNITIES? CONTACT BEN EDWARDS bedwards@kisacoresearch.com