



# Data Platforms Built for the AI Factory

**Bertrand Ounanian**

System Engineer Mgr

VAST Data

**Nouredine Taguelmimt**

Sr Solution Architect

NVIDIA



The Data Platform  
For The AI Era

1000+

Employees Globally

Cash Flow Positive

A Self-Sustaining Growth Engine

\$9.1 Billion

Series E Valuation

Powered By DASE

The First Distributed Systems Architecture Designed for AI



Gartner Peer Insights

84

Net Promoter Score  
(validated by OCX Cognition)



Powering The Most Data-Intensive Enterprise and Cloud Providers

# VAST Data Platform

Multi-Tenant, Zero-Trust, Global Scale



## DataStore

Multi-Protocol File, Object & Block Storage



## DataBase

A Transactional Data Warehouse



## DataEngine

Parallel, Containerized Computing



## DataSpace

Centrally Managed Global Namespace

SFO-AWS



aws

NYC-DC



OnPrem

LON-DC



OnPrem

DUB-C42



core42

TOK-SB



SoftBank



# World's Leading AI Organizations Run on VAST Data

## Model Builders



Powering 100,000+ GPU clusters for the world's leading foundation model builders

## AI CSPs



Providing high-performance data services to 30+ 'neoclouds' all around the world

## Enterprise AI



Real-time Data Platform with **NVIDIA NIMs** for unstructured data in the enterprise

# The AI Era Needs a New Data Stack

## The New AI Data Stack

### Data Activation



Event Broker



Runtime

Real-time processing requires an event-driven architecture to invoke functions for inferencing, embedding, or custom workflows when data is created or altered

### Semantic Layer



Vector Database



KV Cache

Store semantic meaning of unstructured data via embeddings for similarity search and retrieval. New KV caching needed for faster attention computation.

### Data Foundation



Enterprise Storage  
(Multi-Protocol)

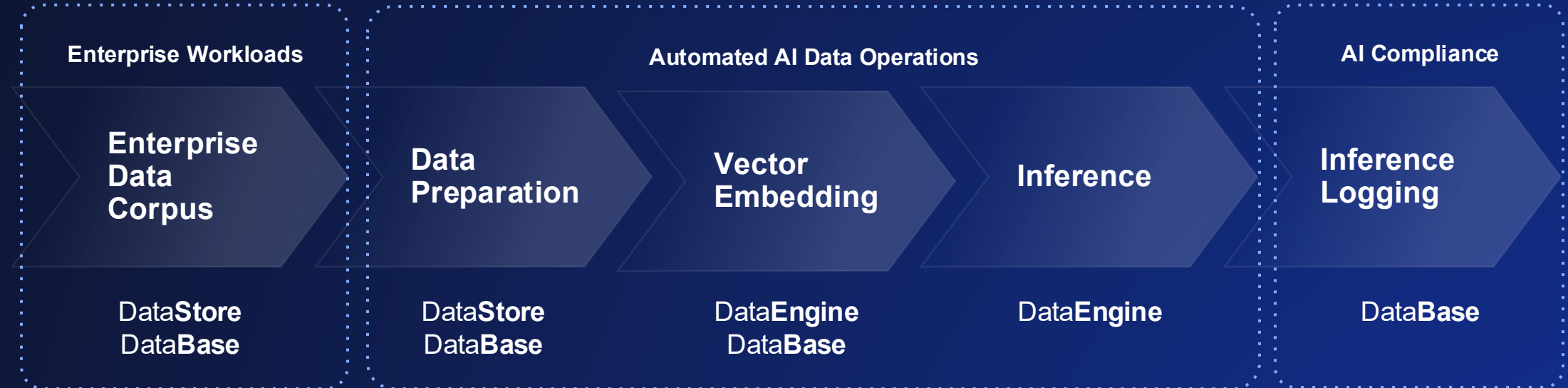


Enterprise Data  
Warehouse

Efficient storing of all data types to power AI learning and reasoning. AI requires fast access to both historical and real-time data and often requires re-indexing multiple times over.

# VAST Data Platform

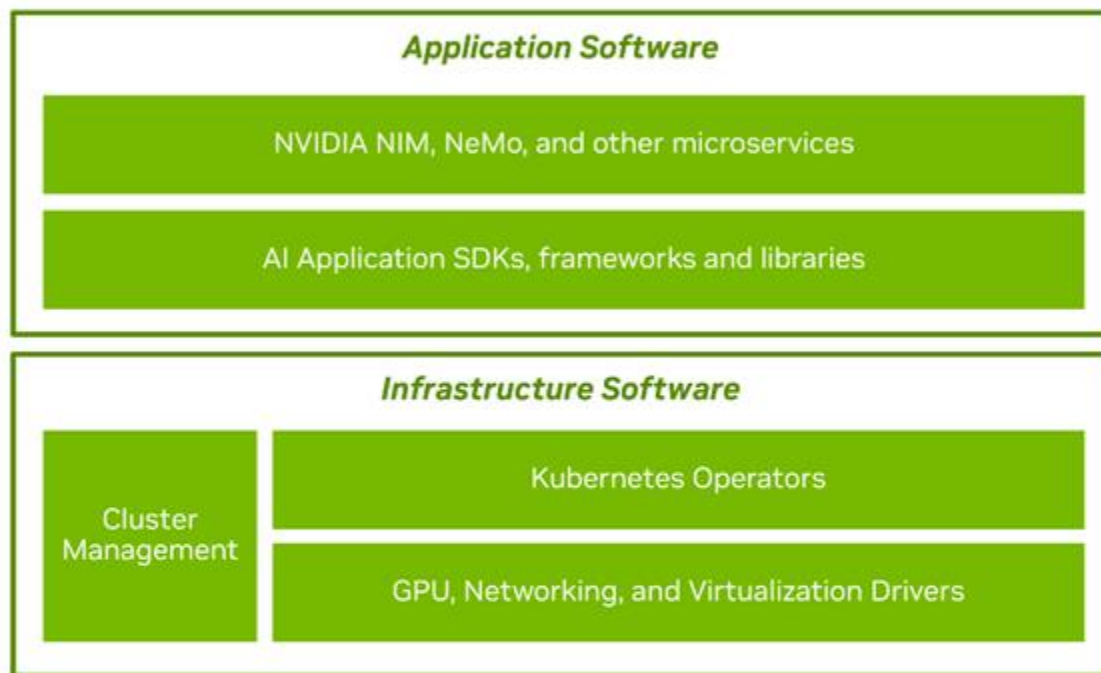
Any Data. Any Workload. No Compromises





# NVIDIA AI Enterprise

What's Inside?



Cloud



Data Center



Edge



Workstation

## Application Software

**Optimized NIM and NeMo microservices** enhance model performance and speed time to deployment for generative AI

**SDKs, Frameworks, and libraries** across many domains, including speech AI, route optimization, cybersecurity, and vision AI

## Infrastructure Software

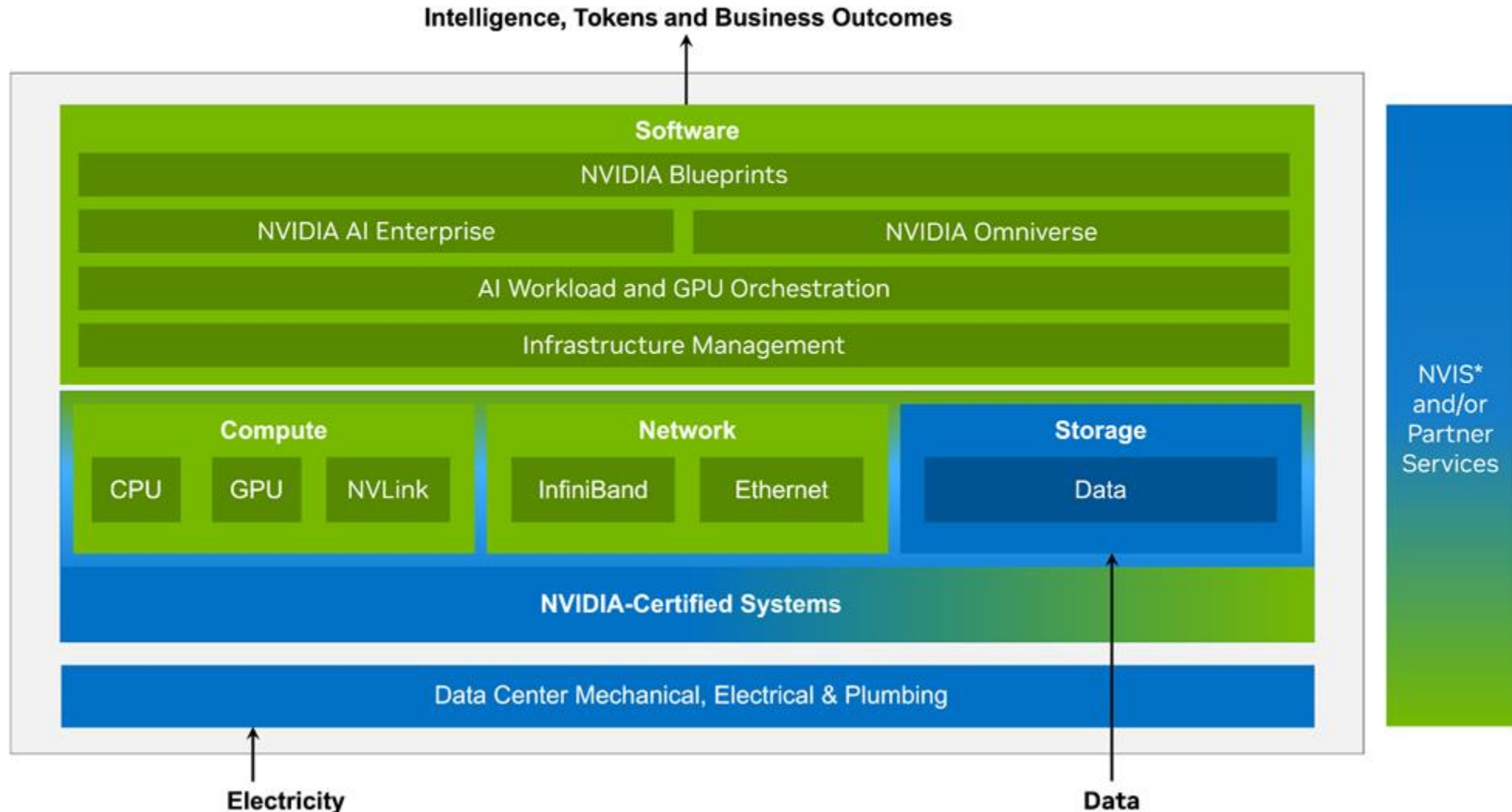
**Kubernetes operators** for managing GPU and networking in containers and the lifecycle of microservices and AI pipelines

**Cluster Management software** to provision and monitor servers at scale

**Drivers** to optimize utilization of NVIDIA GPUs and Networking in bare metal and virtualized environments.

# Full Stack approach for AI Factories

Built on customer-validated data center reference architectures



\*NVIS = NVIDIA Infrastructure Specialists



Link:  
[Jensen Huang keynote showcases the power of InsightEngine through a demo](#)

# VAST InsightEngine

Real-Time Data Streams

Powered With NVIDIA AI Enterprise (NIMS)

- Multi-Modal PDF Extraction
- Video Search & Summarization

AI Apps

AI Apps

AI Apps

AI Apps

AI Apps

Human and Agentic  
Inference Operations

VAST Data Platform

Realtime Structured Data

Atomic Security

DataBase



Tables



Vectors



Graphs

Realtime Unstructured Data

DataStore



Files



Objects

NVIDIA NIM Catalog

Vectors  
Graphs  
Indexes



NIM Semantic  
Engine

chunk-vector/graph-index



Milliseconds

VAST  
DataEngine  
Trigger

# The AI Factory Data Platform



## NVIDIA NIM

Optimized AI inferencing microservice



## NVIDIA Blueprints

Reference AI workflow templates



## NVIDIA NeMo Platform

Build, customize, and deploy generative AI



Enterprise Storage  
(Multi-Protocol)



Enterprise Data  
Warehouse



Vector  
Database



Event  
Broker



Runtime



KV Cache

NVIDIA SPECTRUM-X NETWORKING



## GPUs (H200, B200, etc.)

AI/ML compute acceleration



## CPUs (Grace, etc.)

Orchestration and general-purpose workloads



## DPU (BlueField)

Security, networking, and storage tasks

## VAST InsightEngine with NVIDIA

### Event-Driven Architecture

Invokes NVIDIA NIM in real-time for inferencing and embedding on data creation or modifications

### Unifies End-to-End Security

Enforces IT policies through AI semantic layers automatically

### Scales to Handle All Your Data

Supports billions of files and trillions of embeddings within one platform

### Designed for All Data Types

Store structured and unstructured data via all storage protocols via a single management view

### NVIDIA Certified and Validated

NVIDIA DGX SuperPOD, NVIDIA Cloud Partners, and NVIDIA Enterprise Reference Architectures

