# About Me

**Technology**
- Machine Learning/Deep Learning Hardware and Software Infrastructure
- Digital Twin / Metaverse Hardware and Software
- AI Vertical and Horizontal use case applications

**Location**
- HO: Based in Sicily, Italy
- HQ: Lenovo (Italy) Srl, Via S. Bovio, 3, 20054 Segrate MI

**Professional and Educational Background**
- PhD in Neuroscience and Neurophysiology
- Researcher, Lecturer, Reviewer and Associate Editor in neuroscience and neurophysiology
- Extensive professional experience in Immersive technology applied to pre-clinical research and M&E

**Passion**
- Climbing, Boxing, Trekking, Yoga Nidra
- Avid Book Reader and Movie Watcher
- Photogrammetry, VR Game Dev, Coding
- VR Game Player

## Valerio Rizzo, PhD

| EMEA HEAD of AI & Metaverse SME
| Lenovo (Italy) Srl | ISG
| **email:** vrizzo@lenovo.com

# Unlock the future

# From Fiction to Science



**GREG EGAN**

Beyond space, time, eternity – the ultimate creation dream

**PERMUTATION CITY**



ARE YOU LIVING IN A COMPUTER SIMULATION?

BY NICK BOSTROM

[Published in *Philosophical Quarterly* (2003) Vol. 53, No. 211, pp. 243-255. (First version: 2001)]

This paper argues that *at least one* of the following propositions is true: (1) the human species is very likely to go extinct before reaching a "posthuman" stage; (2) any posthuman civilization is extremely unlikely to run a significant number of simulations of their evolutionary history (or variations thereof); (3) we are almost certainly living in a computer simulation. It follows that the belief that there is a significant chance that we will one day become posthumans who run ancestor-simulations is false, unless we are currently living in a simulation. A number of other consequences of this result are also discussed.

## I. INTRODUCTION

Many works of science fiction as well as some forecasts by serious technologists and futurologists predict that enormous amounts of computing power will be available in the future. Let us suppose for a moment that these predictions are correct. One thing that later generations might do with their super-powerful computers is run detailed simulations of their forebears or of people like their forebears. Because their computers would be so powerful, they could run a great many such simulations. Suppose that these simulated people are conscious (as they would be if the simulations were sufficiently fine-grained and if a certain quite widely accepted position in the philosophy of mind is correct). Then it could be the case that the vast majority of minds like ours do not belong to the original race but rather to people simulated by the advanced descendants of an original race. It is then possible to argue that, if this were the case, we would be rational to think that we are likely among the simulated minds rather than among the original biological ones. Therefore, if we don't think that we are currently living in a computer simulation, we are not entitled to believe that we will have descendants who will run lots of such simulations of their forebears. That is the basic idea. The rest of this paper will spell it out more carefully.

Apart form the interest this thesis may hold for those who are engaged in futuristic speculation, there are also more purely theoretical rewards. The argument provides a stimulus for formulating some methodological and metaphysical questions, and it suggests naturalistic analogies to certain traditional religious conceptions, which some may find amusing or thought-provoking.

The structure of the paper is as follows. First, we formulate an assumption that we need to import from the philosophy of mind in order to get the argument started. Second,



PHOTO: Terry Schneider, Associate Technical Fellow in Boeing Research & Technology, demonstrates computer modeling used to develop new materials at the molecular level. Images on the screen show the molecular structure of resin polymers that bond carbon fibers in composite structures. MARIAN LOCKHART/BOEING

## Atoms to airplanes

New structures technologies, developed across Boeing, are helping accelerate product development **By Bill Seil**

Terry Schneider, an Associate Technical Fellow in Boeing Research & Technology, works in "atoms to airplanes" modeling, or the complete process of modeling an airplane computationally from a molecular level up to the full-scale, complete airframe.

One important goal of this work is to optimize the chemistry of polymers to increase the load-carrying capability of the carbon fiber in composites, which could significantly reduce the weight of next-generation composite structures.

"This is exciting work because we're able to rapidly assess hundreds of polymer candidates in a matter of weeks—a process that might take years in a lab," Schneider said. "We're also able to quickly determine their performance in large-scale laminated structures and screen for the best-performing candidates. This opens the door to huge cost savings in the future."

Work such as this demonstrates the benefits to Boeing generated by the company's enterprisewide approach to making

research investments in key areas such as structures, a term that describes the physical airframe components of airplanes and other aerospace products. Critical aviation design issues—including weight, reliability and safety—all depend on the quality of research and planning that drives structures engineering.

Boeing has long been a leader in structures technology, and research conducted throughout the enterprise has steadily improved the design of structures and the materials used to make them. The challenge today is to increase the company's competitive edge by investing in research that generates maximum benefit for Boeing's range of products, both commercial and military.

That's why, in 2008, the company created its Enterprise Technology Strategy (ETS), which takes a coordinated, "One Company" approach to technology development. The strategy is built around eight technology areas, or domains, that support Boeing's many business programs and can create a sustainable technical competitive advantage that helps the company grow.

# From Digital Twin to Industrial Metaverse



Industrial Metaverse

Whole-System DT

Immersive DT

Digital Twin

"

*Connected whole-system digital twin with functionalities to interact with the real system in its environment, allowing decision makers to better understand the past and forecast the future."*

Arthur D. Little

## SIEMENS

A virtual world in which we can **interact in real time with photorealistic, physics-based digital twins** of our real world. We believe **digital twins** are **the building blocks for the Metaverse.**

## NVIDIA

Industrial Metaverse enables industrial companies of all sizes to create **closed-loop digital twins with real-time performance data, ideal for running simulations** and AI-accelerated processes for advanced applications such as **autonomous factories that rely on intelligent sensors and connected devices.**

## IndustrialMetaverse.org

A real-time, persistent simulation space that is the **sum of all virtual worlds, digital twins, and augmented reality** that **connects digital economic assets and infrastructure** on a global scale in the **industrial and commercial setting.**

## Microsoft

Industrial Metaverse enables **humans and AI to work together to design, build, operate, and optimize physical systems** using digital technologies.

## sme

A **systematic discipline that combines hardware [...] data conversions** through analytics/machine learning, **time histories** through cyber-infrastructure, **cognition** through human-machine interface, **and configuration** through the Metaverse.

## COSMO TECH

The Industrial Metaverse enables the creation of **digital twins of places, processes, real-world objects, and** the **humans** who interact with them.

Source: Arthur D. Little

"

*A massively <u>scaled</u> and <u>interoperable</u> network of real-time rendered 3d virtual worlds that can be experienced <u>synchronously and persistently</u> by an effectively <u>unlimited number of users</u> with an <u>individual sense of presence</u> and with <u>continuity of data</u>, such as identity, history, entitlements, objects , communications and payments* "

*Matthew Ball, The Metaverse*

# Anatomy of the Metaverse



Services, Games, Shopping, Events, more → Experience

Ad Network, Socials, Rating, Stores, Agents → Discover

Design and edit tools, Assets Markets, Platforms → Creator Economy

Spatial Computing → Game Engines, Multitasking UI, Geospatial Coherence, AR/VR/MR

IoT, Microservice, Blockchain, NFTs → Decentralization

User Interface → Mobile, BCI, Haptic, Voice, Gesture

5G/6G, WiFi 6, Cloud, 7nm to 1.4 nm, XPUs, Edge Computing, Storage → Infrastructure

Thien Huynh-The et al. "Artificial Intelligence for the Metaverse: A Survey" arXiv:2202.10336v1 [cs.CY] 15 Feb 2022

# Metaverse System Model

# The Intertwined Nature of Metaverse and AI

**Sensors/IoT/Sim**

**Uses**

**Digital Twin/Metaverse**

**Generate**

**Big Data**



**Creates/ Operates/ Monitors**

**Employs**

**Validates**

**Generates**

**AI / GenAI**

# AI value for the Metaverse

Personalization
Education
Assistance

Recommendation

Content Generation

Compute Speed-up

Smart Contracts

Inclusivity&Intuitivity

AI ops, MLops

Experience

Discover

Creator Economy

Spatial Computing

Decentralization

User Interface

Infrastructure

Services, Games,
Shopping, Events, more

Ad Network, Socials,
Rating, Stores, Agents

Design and edit tools,
Assets Markets,
Platforms

Game Engines,
Multitasking UI, Geospatial
Coherence, AR/VR/MR

IoT, Microservice,
Blockchain, NFTs

Mobile, BCI, Haptic, Voice,
Gesture

5G/6G, WiFi 6, Cloud, 7nm to 1.4
nm, XPUs, Edge Computing

Thien Huynh-The et al. "Artificial Intelligence for the Metaverse: A Survey" arXiv:2202.10336v1 [cs.CY] 15 Feb 2022

# How Today's AI is Shaping Tomorrow's Possibilities

**3D Modeling & Visualization**

**Decentralized Computing**

**Network Optimization**

**Confidential AI Solutions**



**INDUSTRIAL METAVERSE**

**Spatial Computing**

**Human-Machine Interactivity**

**Physically Accurate Simulations**

**Realistic Interactive Virtual Entities**

# How Today's AI is Shaping Tomorrow's Possibilities

**3D Modeling & Visualization**

**Decentralized Computing**

**Network Optimization**

**Confidential AI Solutions**

**INDUSTRIAL METAVERSE**

**Spatial Computing**

**Human-Machine Interactivity**

**Physically Accurate Simulations**

**Realistic Interactive Virtual Entities**

# Advances in Neural Rendering and Mesh Generation

## PixelNeRF (2021)



## Instant Ngp (2022)



## Neuralangelo (2023)



## Magic3D (2023)



| Metric | PixelNeRF | Instant NGP | Neuralangelo |
|---|---|---|---|
| Rendering Time (ms) | 10-30 per pixel | <1 per pixel | ~100-500 per pixel |
| Scene Complexity | High | Medium-High | Very High |
| Photorealism | No/ Limited | Yes | Yes |
| Real-time Capability | No | Yes | No |

A. Yu, V. Ye, M. Tancik and A. Kanazawa, "pixelNeRF: Neural Radiance Fields from One or Few Images," 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Nashville, TN, USA, 2021, pp. 4576-4585, doi: 10.1109/CVPR46437.2021.00455.
keywords: {Convolutional codes;Solid modeling;Computer vision;Three-dimensional displays;Image resolution;Computer architecture;Benchmark testing},
"Instant Neural Graphics Primitives with a Multiresolution Hash Encoding" Thomas Müller et al. ACM Transactions on Graphics (SIGGRAPH), July 2022a
Li, Zhaoshuo & Müller, Thomas & Evans, Alex & Taylor, Russell & Unberath, Mathias & Liu, Ming-Yu & Lin, Chen-Hsuan. (2023). Neuralangelo: High-Fidelity Neural Surface Reconstruction.
Lin, Chen-Hsuan & Gao, Jun & Tang, Luming & Takikawa, Towaki & Zeng, Xiaohui & Huang, Xun & Kreis, Karsten & Fidler, Sanja & Liu, Ming-Yu & Lin, Tsung-Yi. (2022). Magic3D: High-Resolution Text-to-3D Content Creation. 10.48550/arXiv.2211.10440.

# Towards Real-Time Physically Accurate Simulations


Approximating mesh-motion Laplacian mesh motion solver in OpenFOAM with MLP


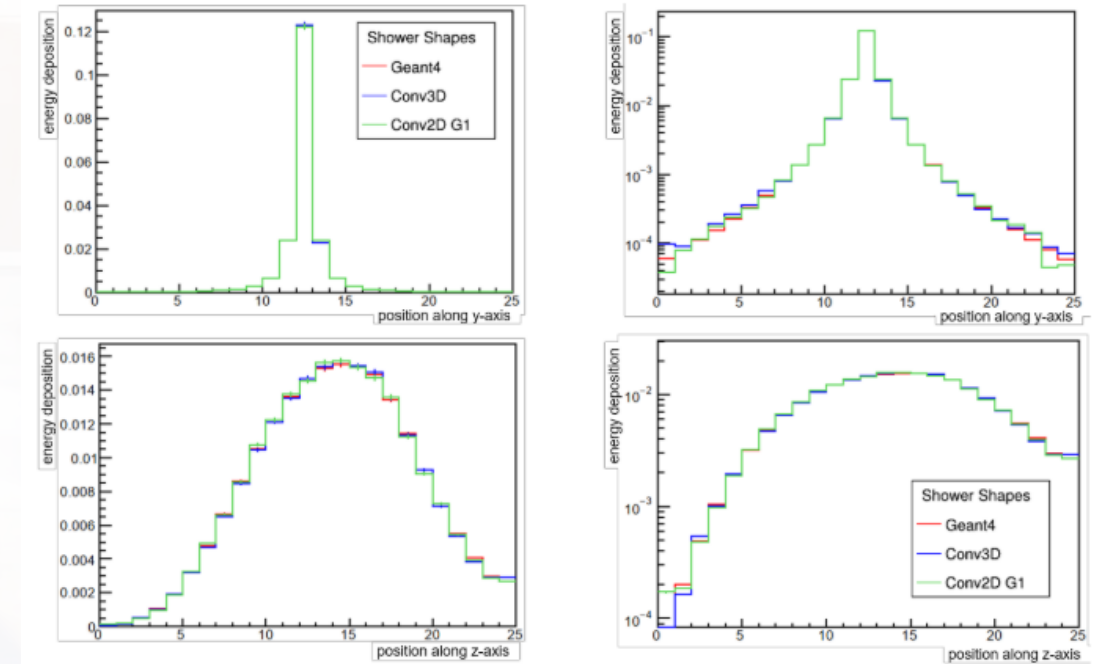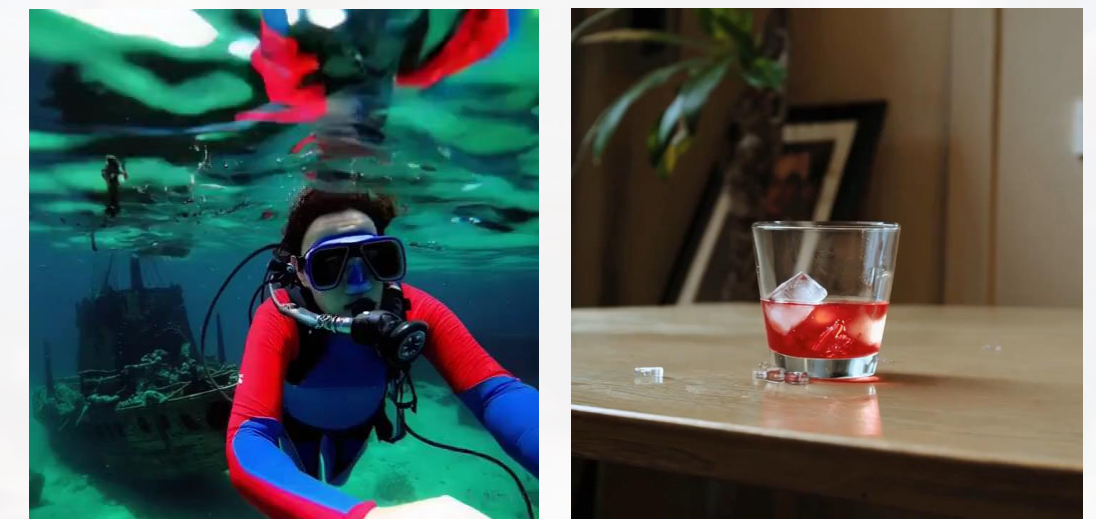Using CNN to Solve Euler-Lagrange, Momentum Transfer, and Incompressible RANS Equations

Gen AI

AI / ML

16


Simulating high energy physics calorimeter detector outputs with 2D GAN


Video generation models as general purpose simulators of the physical world?

Maric, T., Fadeli, M.E., Rigazzi, A. et al. Combining machine learning with computational fluid dynamics using OpenFOAM and SmartSim. Meccanica (2024). https://doi.org/10.1007/s11012-024-01797-z
Rojek, K., Wyrzykowski, R., Gepner, P. (2021). AI-Accelerated CFD Simulation Based on OpenFOAM and CPU/GPU Computing. In: Paszynski, M., Kranzlmüller, D., Krzhizhanovskaya, V.V., Dongarra, J.J., Sloot, P.M.A. (eds) Computational Science – ICCS 2021. Lecture Notes in Computer Science(), vol 12743. Springer, Cham. https://doi.org/10.1007/978-3-030-77964-1_29
Rehm, Florian et al. "Validation of Deep Convolutional Generative Adversarial Networks for High Energy Physics Calorimeter Simulations." AAAI Spring Symposium: MLPS (2021).
https://openai.com/index/video-generation-models-as-world-simulators/

# Tensors Reshape Compute Architectures



NVIDIA

AMD

Intel

NVIDIA

Intel

AMD

NVIDIA

# AI-optimized Portfolio from Model Development to Inferencing
## *80+ new and enhanced Infrastructure platforms – Pocket to Cloud, Edge to Core*
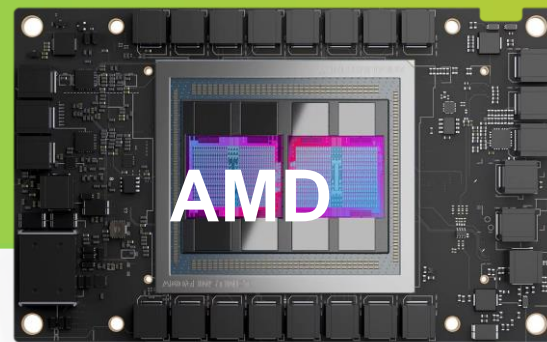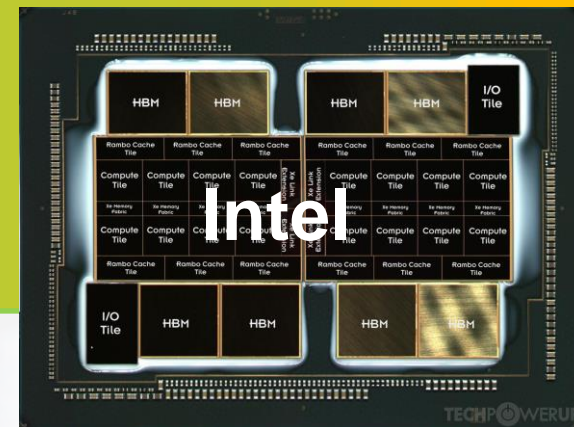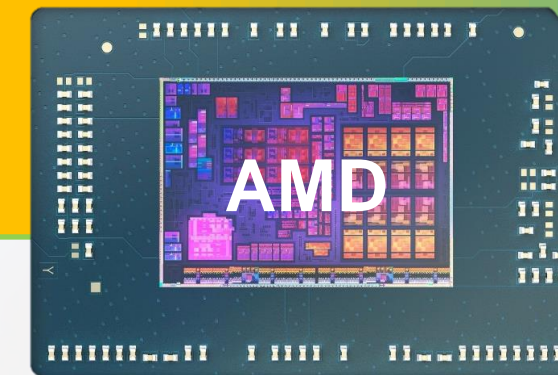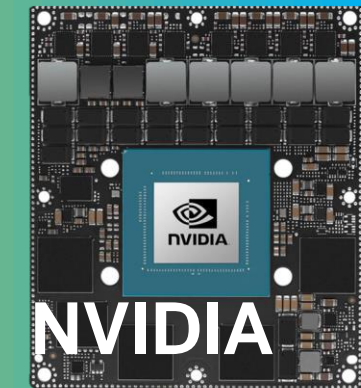
### Data Management

**Solutions**
High Performance File System (w/WEKA)
Object Storage Solutions (w/Cloudian)
DSS-G / Spectrum Scale
BeeGFS

**DM & DE**         DE6600
DM7100F         DE6600
DM5100
DG7000

### ML & Data Analytics

**4-socket**
SR850 V2
SR850 V3 Intel
SR860 V2
SR860 V3 Intel

**2-Socket**
SR650 V2
SR650 V3 Intel
SR655
SR655 V3 AMD
SR665
SR665 V3 AMD

**ThinkSystem**

### Deep Learning Training HGX

ST650 V3 Intel
SR670/75 V3 4-8x PCIe
SR670/75 V3 4-GPU HGX

NEW SR680a V3 8-GPU HGX
NEW SR685 V3 8-GPU HGX

SD650-I V3
SD665-N V3

**NEW** SR780a 8-GPU HGX

### Data Science Workstation

**Edge**         P1 Gen5
NEW P3 Ultra   P1 Gen6
NEW P3 Tiny
**Desktop**
NEW PX
NEW P7
P620
NEW P5
NEW P3 Tower

### ThinkPad with Neural Processing Units

ThinkPad X13s Gen1 – 15 TOPS
ThinkPad Z13 Gen2 – 11 TOPS
ThinkPad Z16 Gen2 – 11 TOPS
ThinkPad T14s AMD Gen4 – 11 TOPS
ThinkPad T14 AMD Gen4 – 11 TOPS
ThinkPad T16 AMD Gen2 – 11 TOPS
ThinkPad X13 AMD Gen4 – 11 TOPS

### Edge AI

**Server**          SE70 AWS
SE350             Panorama
NEW SE350 V2
NEW SE360 V2
SE450
NEW SE455

**Clients**
SE10, SE10-I
M90
SE30
SE50
SE70

**AI Appliance**

### Appliances

**ThinkAgile**
**MX Systems**
**(Microsoft)**
MX3330-F
MX3330-H
MX3331-F
MX3331-H
MX3530-F
MX3530-H
MX3531-F
MX3531-H

**ThinkAgile**
**HX Systems**
**(Nutanix)**
HX1330
HX1331
HX2330
HX2331
HX3330
HX3331
HX5530
HX5531

**ThinkAgile**
**VX Systems**
**(VMware)**
VX3331
VX3530-G
VX7531

Lenovo
**ThinkStation**

Lenovo
**ThinkPad**
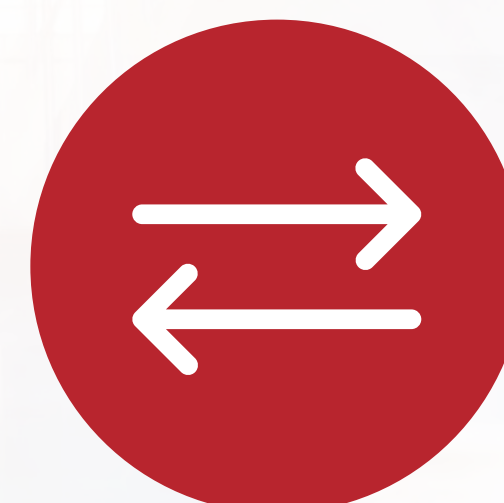
Lenovo
**ThinkEdge**

**ThinkAgile**

# Challenges ahead

Scalability
& Energy Efficiency

Security &
Privacy

Interoperability &
Standards

Compute & Storage
Optimization

Ethic & Regulations
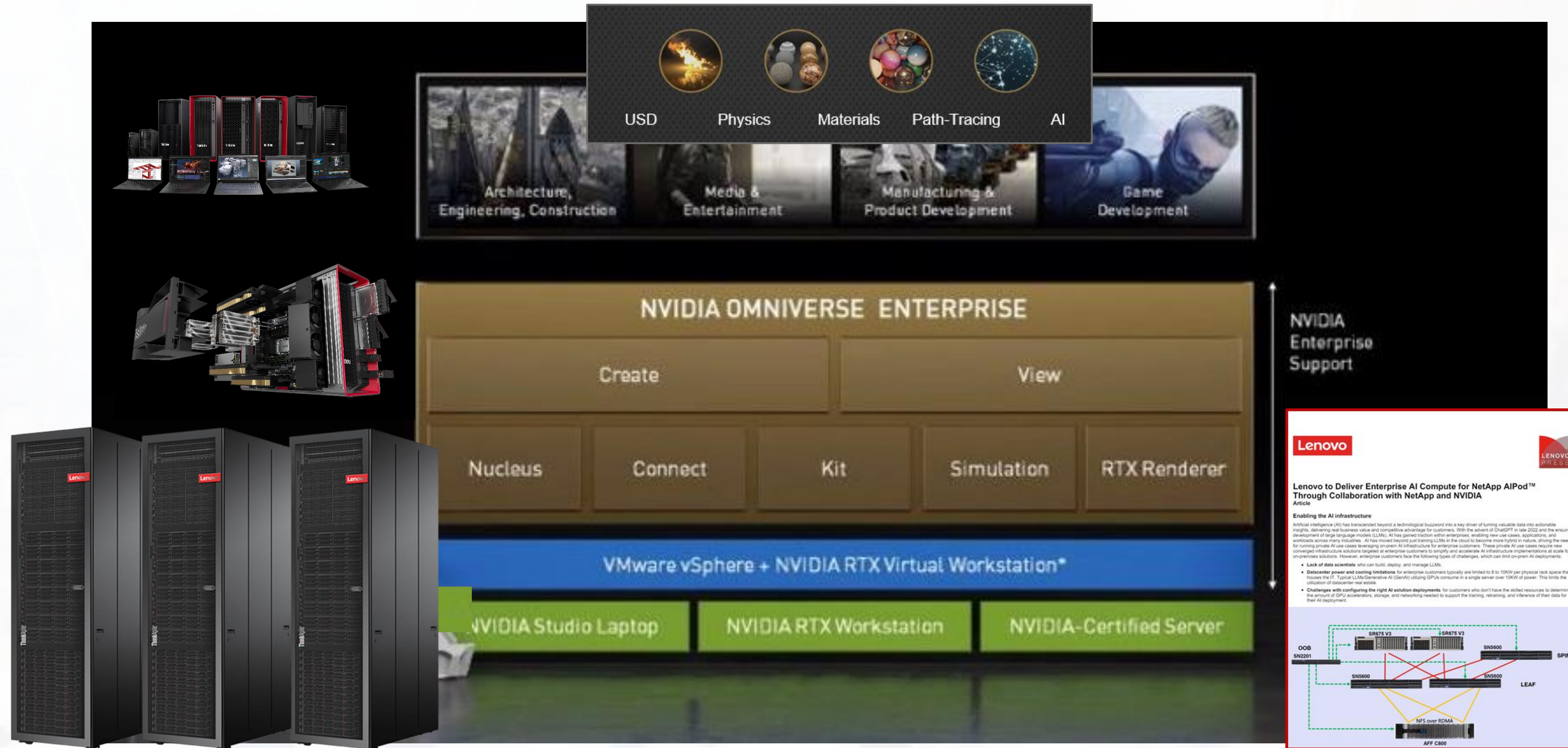
# Lenovo E2E – OVX Infrastructure Solutions
## *Through Collaboration with NetApp and NVIDIA*

# The benefits of MV tech application embrace all industries.

| Automotive | Energy | Infrastructure | Retail | Science |
|---|---|---|---|---|



**Automotive**
- Fast-Track Industrial Factory Planning
- Developing Custom Applications for Factory Planners

**Energy**
- Accelerating Fusion Reactor Design and Development
- Reducing Downtime and Unplanned Maintenance
- Optimizing Wind Farm Design and Electricity Generation

**Infrastructure**
- Transforming Telco Network Planning and Operations
- Simulating and Optimizing Autonomous Railway Networks
- Testing and Optimizing 5G Deployment

**Retail**
- Autonomous Warehouse Robots
- Retail Layout
- Optimizing Distribution Center Throughput

**Science**
- Accelerating Carbon Capture and Storage
- Visualizing High-Resolution, Global-Scale Climate Data
- Accelerating Climate Research
- Visualizing Molecular Dynamics
- Brain Digital Twin

21

# Industrial Metaverse

# Are we there yet?

## Takeaways:

| | |
|---|---|
| **Evolving DT Concept** | The extended and enhanced use of digital twins is at the core of the Industrial Metaverse. AI applications can speed up 3D asset creation and prototyping while providing more intelligent capabilities to DT |
| **AI-Powered Metaverse** | Integrating AI into the HPC framework for the Industrial Metaverse unlocks new capabilities, driving innovation and efficiency in high-fidelity rendering and physical simulations. |
| **Metaverse-Ready Infrastructure** | The key technologies for achieving extended whole-system digital twins are not yet mature, but advances in AI, edge computing, and cloud infrastructure are rapidly closing the gap. |
| **Challenges** | Key issues include security, scalability, latency, costs, skill gaps, and regulatory compliance (including AI and data governance) |
| **Future Trends** | Accelerators mem bw will keep increasing, AI eats HPC, Raytracing engine will be integrated into AI superchips (i.e.: NVIDIA DGX) or Viz card will start employing DGX-like architectures |

22

Forum
TERATEC **24**

# thanks.