
NVIDIA

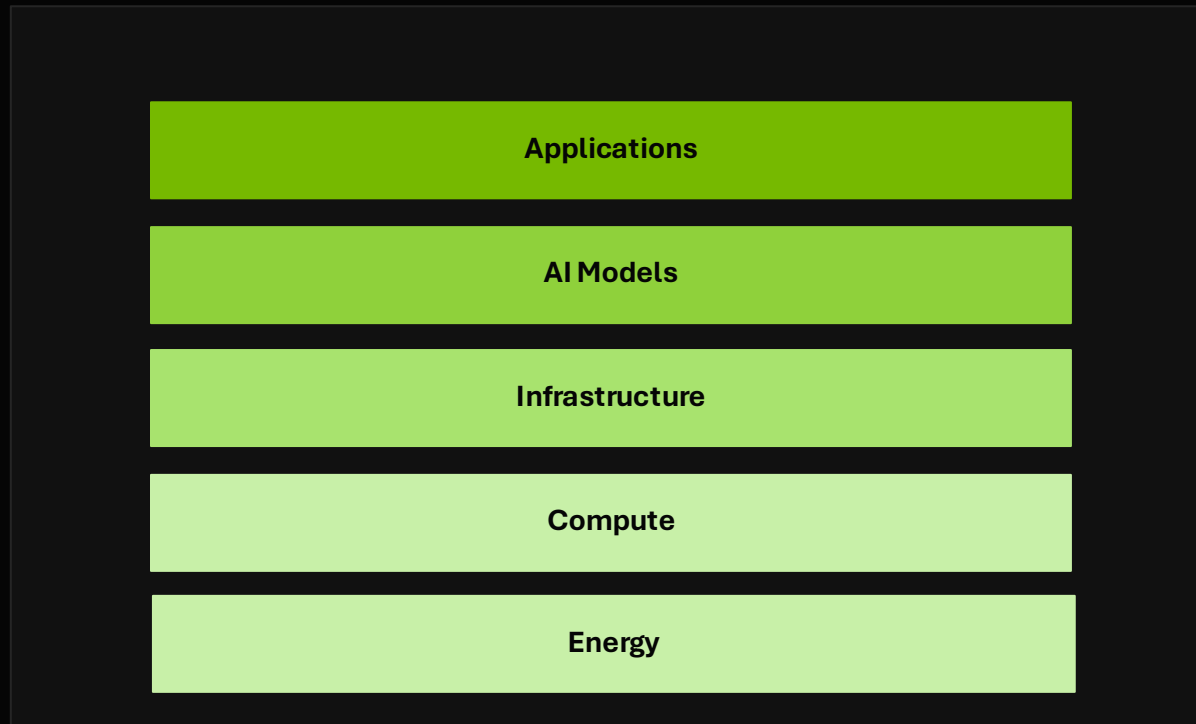
Increasing Compute Density, Agentic and Physical AI

Flexible AI factories for resiliency, affordability, and speed to power

ERCOT Board

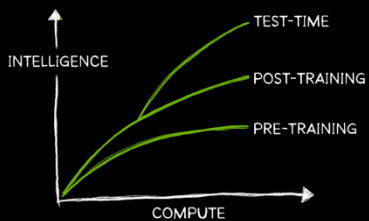
AI is Five Layer Cake

Energy is no more a commodity- it is the foundation of the AI technology stack.



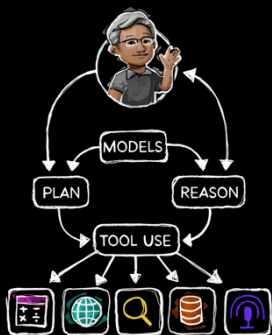
AI Scales Beyond LLMs

Growth in Reasoning model, and Physical AI supercharging Inferencing



3 SCALING LAWS

COMPUTE IS DATA



AI BECOMES AGENTIC



PHYSICAL AI TAKES LEAP



AI LEARNS LAWS OF NATURE



OPEN MODELS REACH FRONTIER

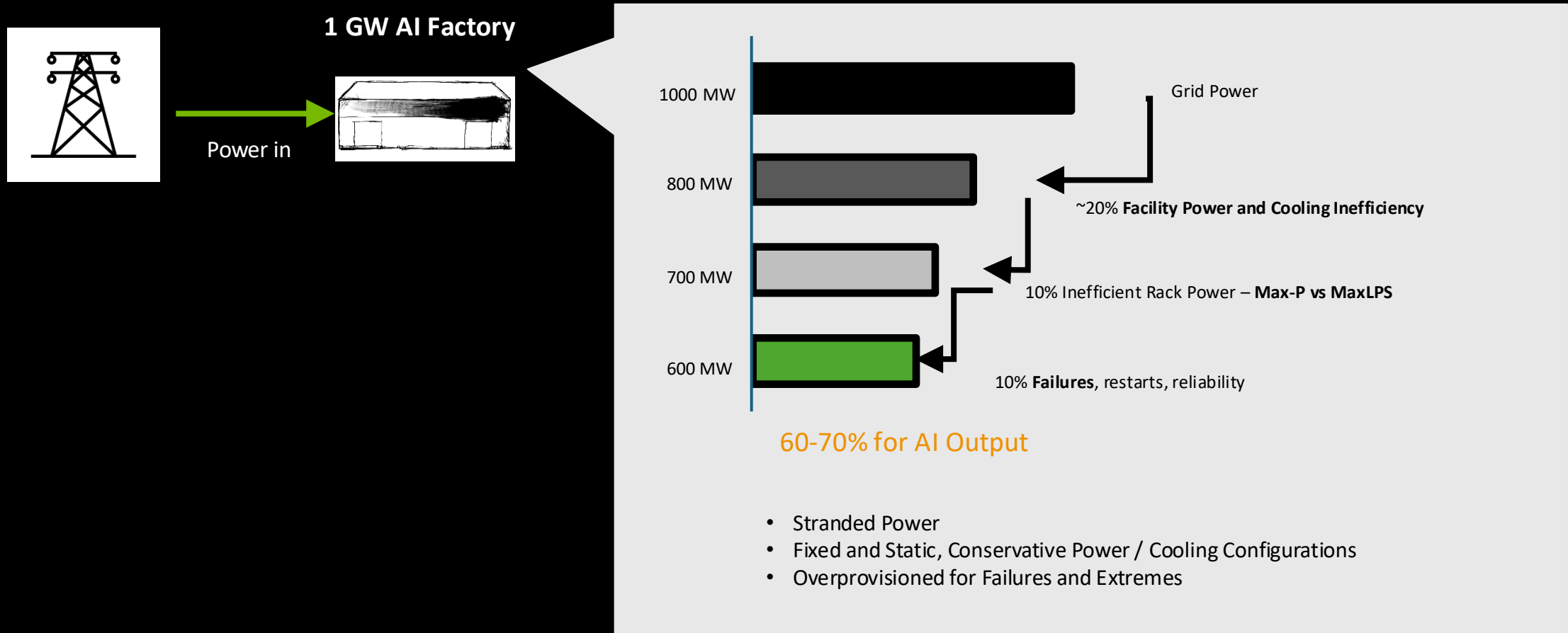
AI Growth Is Not a Bubble

ChatGPT to Claude Code (AI Agents)- 10,000X Inference Compute



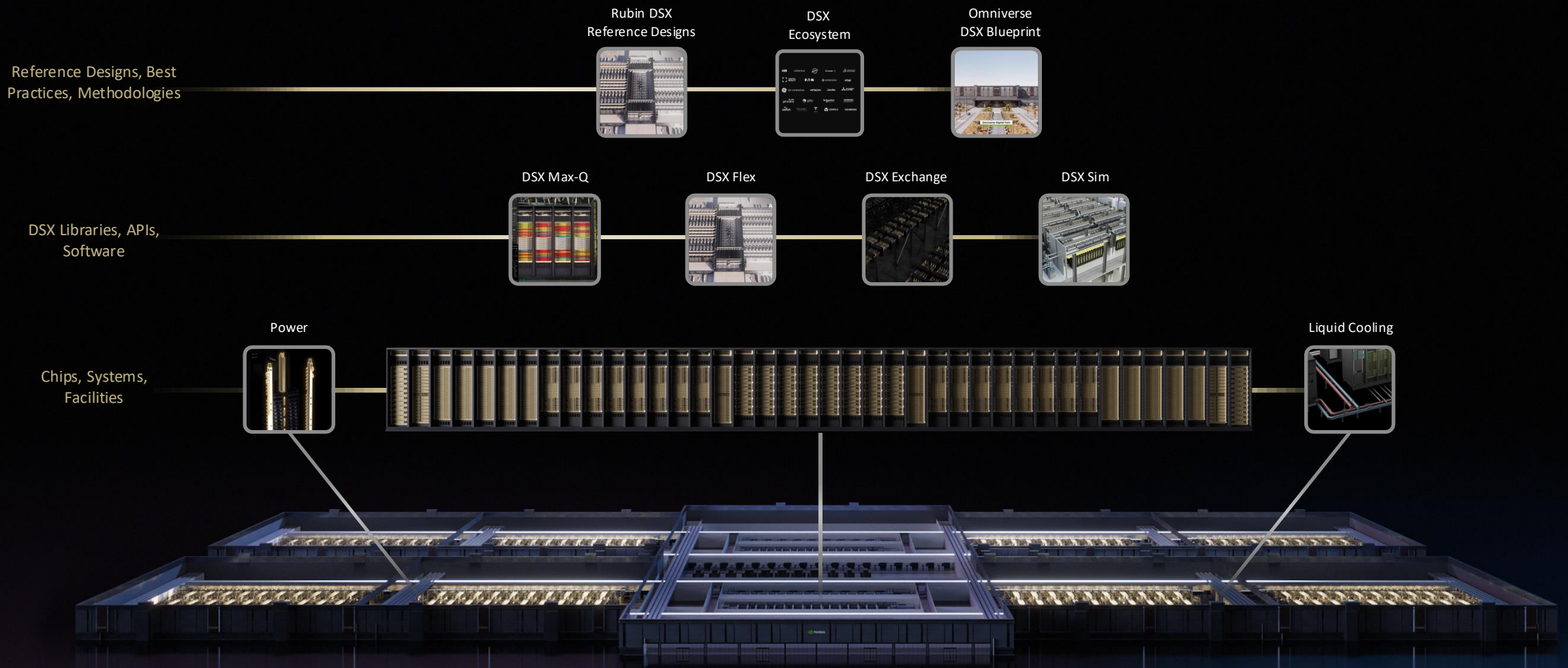
Challenge is to maximize AI tokens per watt of available energy

Facility, rack, downtime inefficiencies cause only a fraction of grid power to contribute to AI output



NVIDIA DSX – Open Reference Architecture for GW AI Factory

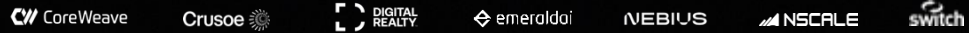
Bringing over 100 Ecosystem Partners adopt standard design to maximize efficiency



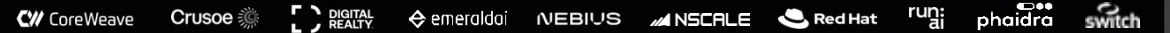
NVIDIA DSX AI Factory Platform

Accelerates Scalable, Energy-Efficient AI Factory Deployment

DSX Flex



DSX Max-Q



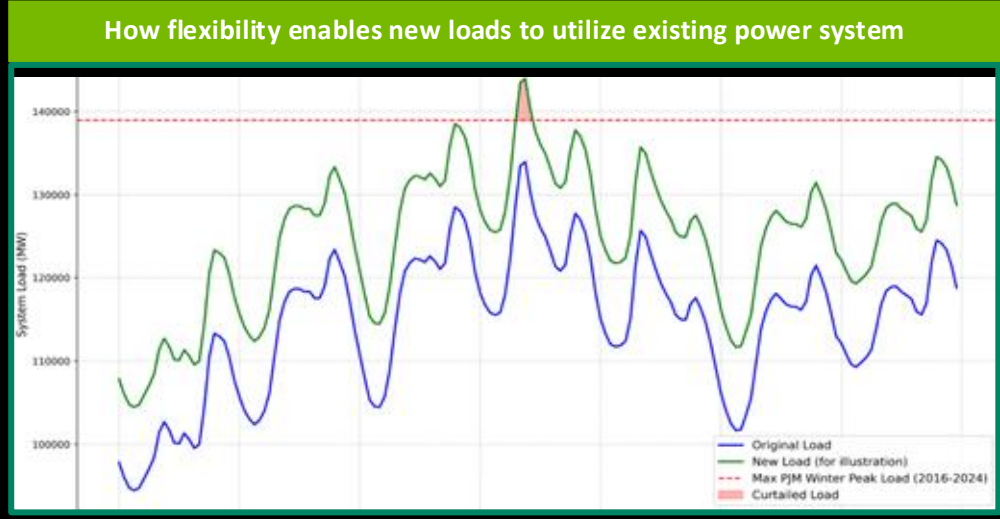
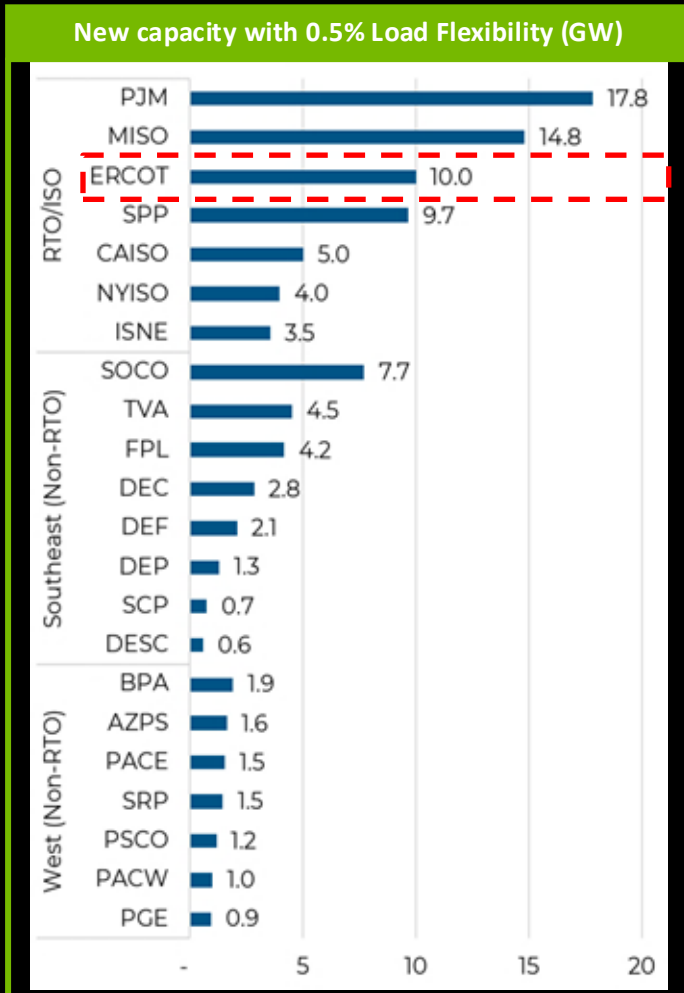
DSX Sim



DSX Exchange



Flexible AI Factories can unlock 10GW or more in ERCOT alone to power the AI revolution, protect affordability, and best utilize the existing power system



First Mover Technology Companies (Duke, 2025)

Category	Examples
Operational flexibility	<ul style="list-style-type: none"> Google deployed a “carbon-aware” temporal workload-shifting algorithm and is now seeking to develop geographic distribution capabilities (Radovanović 2020). Google data centers have participated in demand response by reducing non-urgent compute tasks during grid stress events in Oregon, Nebraska, the US Southeast, Europe, and Taiwan (Mehra and Hasegawa 2023).
	<ul style="list-style-type: none"> Startup companies like Emerald AI are developing software to enable large-scale demand response from data centers through recent advances in computational resource management to precisely deliver grid services while preserving acceptable quality of service for compute users

Source: Norris et al., Duke University, 2025

EPRI DCFlex- Grid-integrated Flexible Data Centers

Developers



Hyperscalers



IPP's



ISO/RTO



Technology Providers



Advisory & Finance



Engineering & Construction

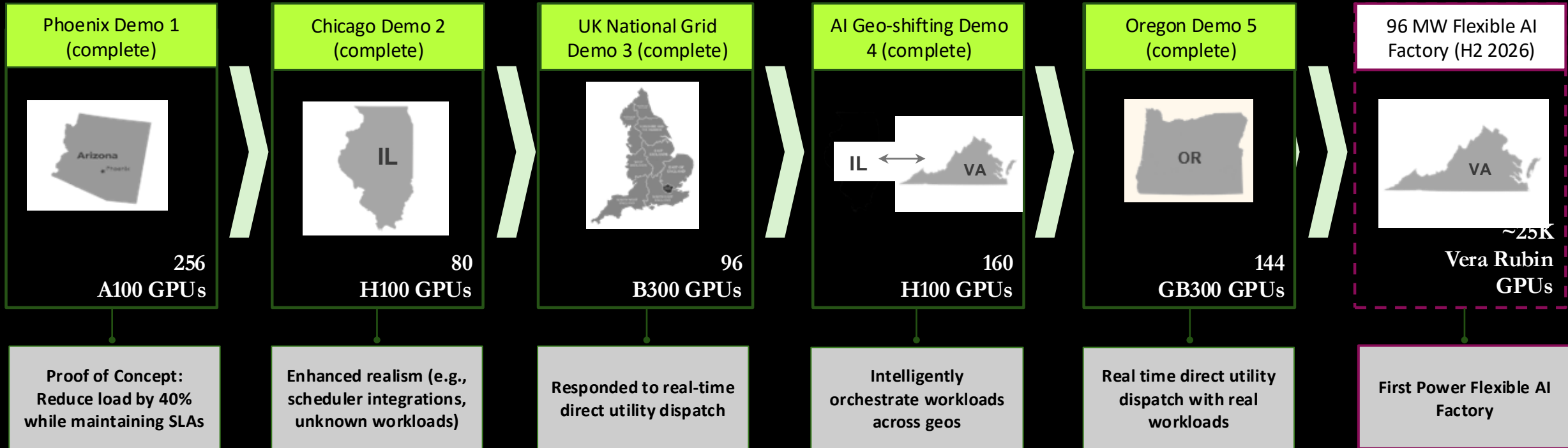


Utilities



DCFlex- NVIDIA and Emerald AI have completed a series of successful demos and deployments ahead of a 96 MW commercial deployment in Virginia

Demo Schedule



Each subsequent Demo has demonstrated our growing capabilities under more challenging circumstances

NVIDIA and Emerald AI have demonstrated AI Factory flexibility globally, including by relieving grid strain at halftime of a football match

Demonstration Goal

Demonstrate that an AI cluster can meet a designated power reduction in response to high periods of electricity use during major TV breaks in the UK

Results

Successfully reduced cluster power by **35% for 30 minutes surrounding halftime** and for **25% for one hour at the end** of a football match, helping combat the 'TV Pickup' phenomenon as viewers fire up electric tea kettles during breaks in TV viewing

Helped ease burden from a National Grid ESO-wide spike of **1 GW at halftime** and **1.6 GW at full-time**

Cluster Power Trace during soccer match (35% and 25% Reductions)



Open Power AI Consortium- EPRI

First Open-Source AI Models for the Power Sector

Domain-Specific Solutions
Improve Operations, Energy Efficiency, and Grid Resilience

Domain-Specific Model
Developed by EPRI, NVIDIA, and Articul8

Domain-Specific Data
Proprietary EPRI Energy and Electrical Engineering Data

~10,000
EPRI Files

400,000+
Images

~230,000
Tables

Texas leading the AI buildout

Building Grid Friendly Flexible AI Factories/Data Centers



Supercharge AI Growth

Reasoning Models, Agents, and Physical AI



Grid Resiliency

Flexible AI Factory, providing Peak load management



Affordability

Increase utilization of existing grid assets