

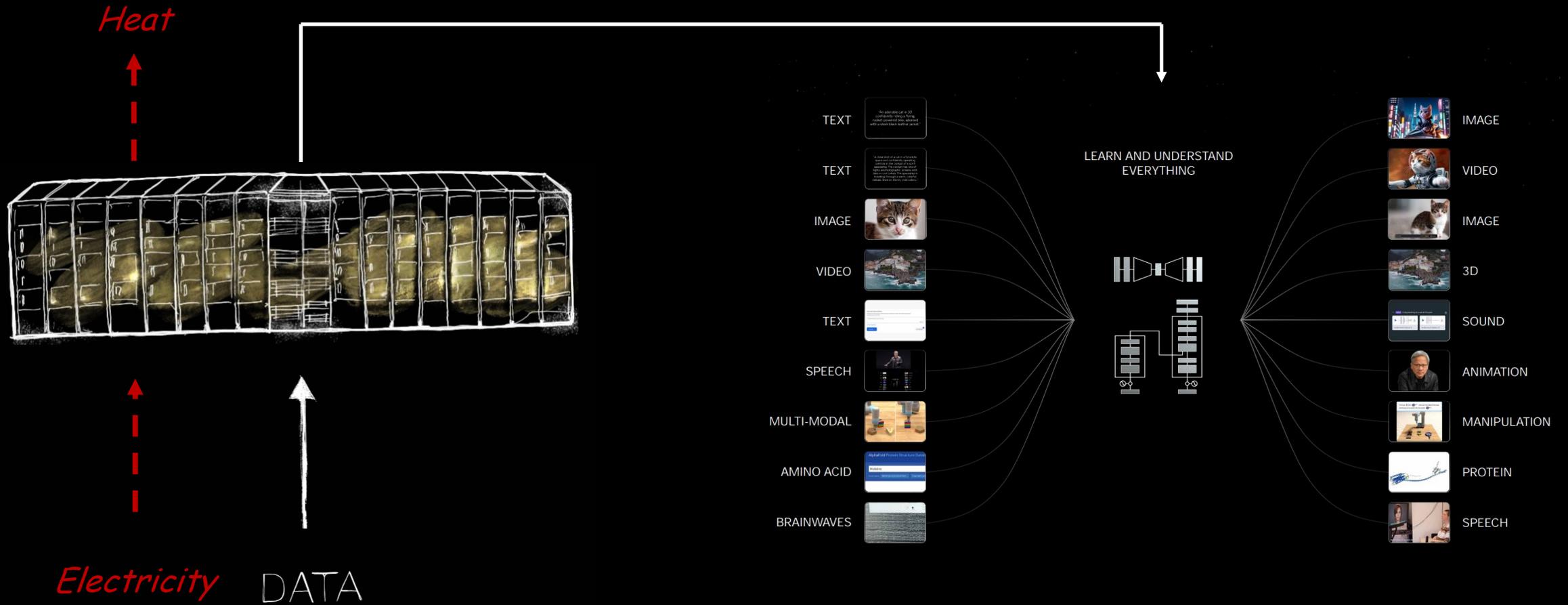
성공적인 엔터프라이즈 생성형 AI 여정을 위한 풀스택 인프라

—
김진효 수석부장, 메이머스트

2024

AI 는 Data, Algorithms, 그리고 AI 인프라를 필요로 합니다

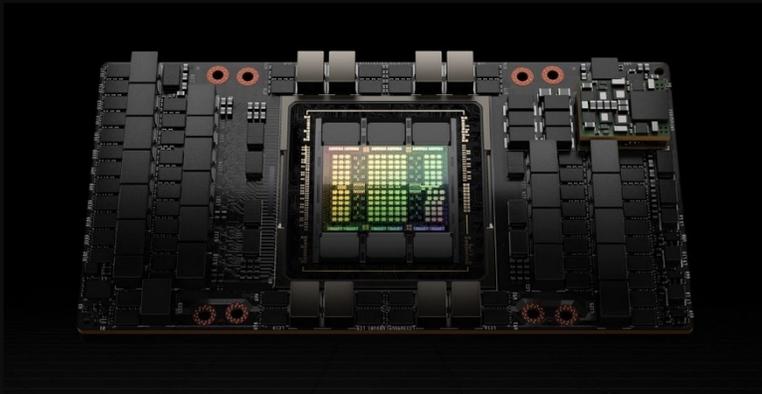
...



• ATTENTION

Industry Dynamics – AI 인프라의 핵심은 GPU

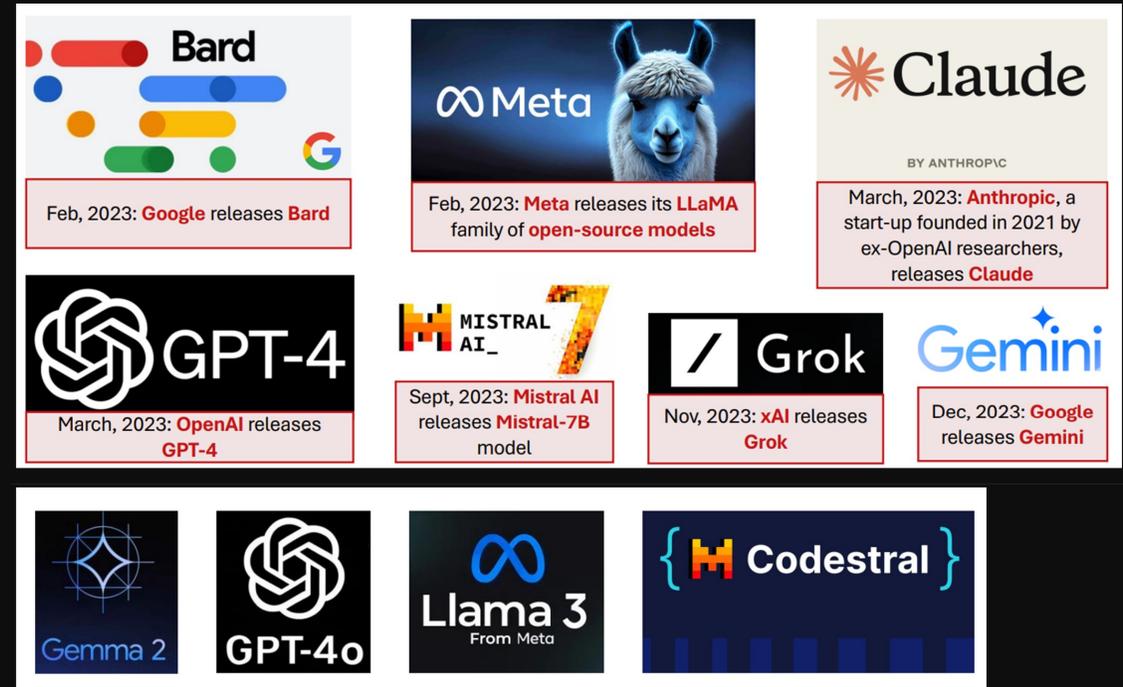
AI 반도체, 네트워크, SW 생태계를 선점한
NVIDIA 지배력은 당분간 유지 전망



2023



2024

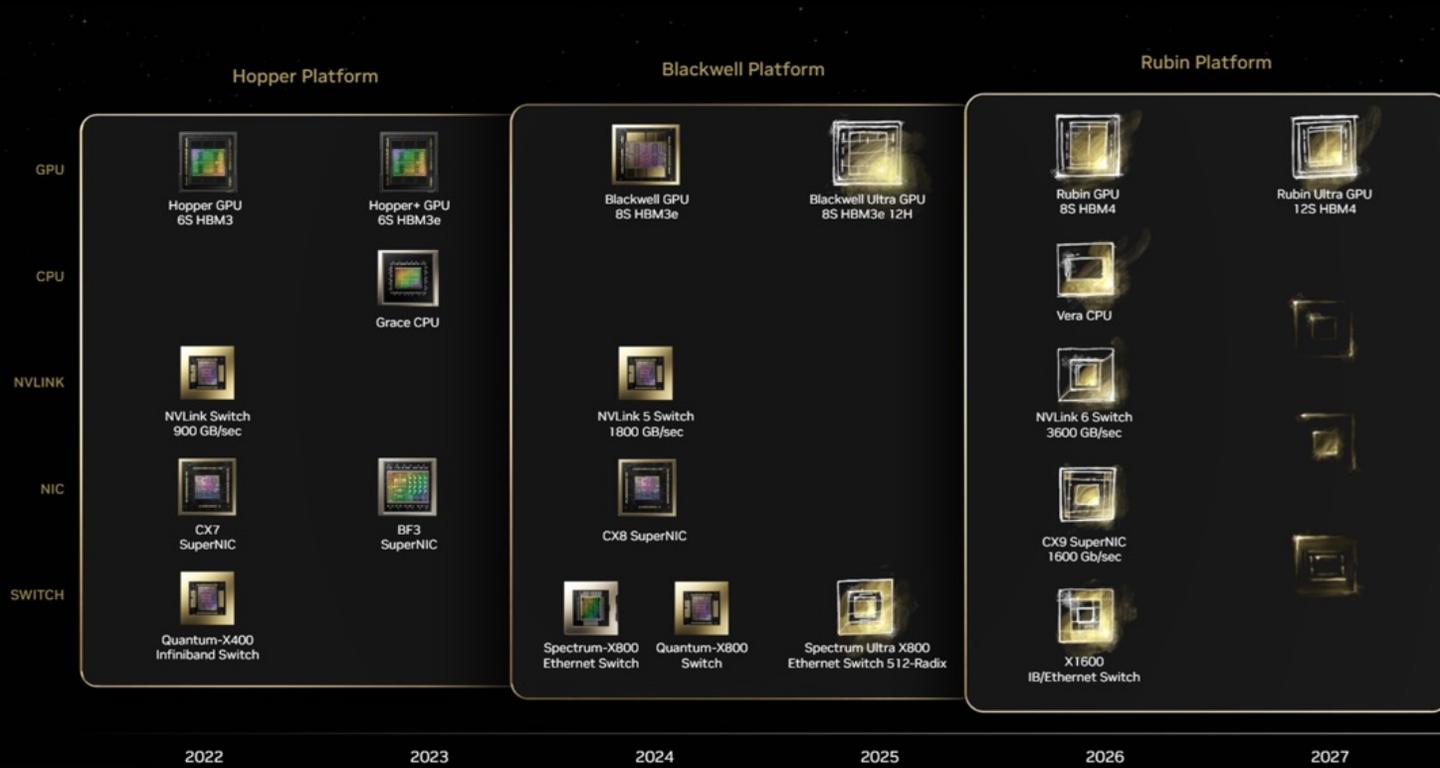


Timeline of AI model releases:

- Feb, 2023: Google releases Bard
- Feb, 2023: Meta releases its LLaMA family of open-source models
- March, 2023: Anthropic, a start-up founded in 2021 by ex-OpenAI researchers, releases Claude
- March, 2023: OpenAI releases GPT-4
- Sept, 2023: Mistral AI releases Mistral-7B model
- Nov, 2023: xAI releases Grok
- Dec, 2023: Google releases Gemini
- 2024: Gemma 2
- 2024: GPT-4o
- 2024: Llama 3 From Meta
- 2024: Codestral

Content Credit: https://lcs2-iiid.github.io/ELL881-AIL821-2401/static_files/presentations/101.pdf

AI Infrastructure is Evolving – (너무) 빠른 진화



DATACENTER SCALE • ONE-YEAR RHYTHM • TECHNOLOGY LIMITS • ONE ARCHITECTURE

Images source – NVIDIA presentations

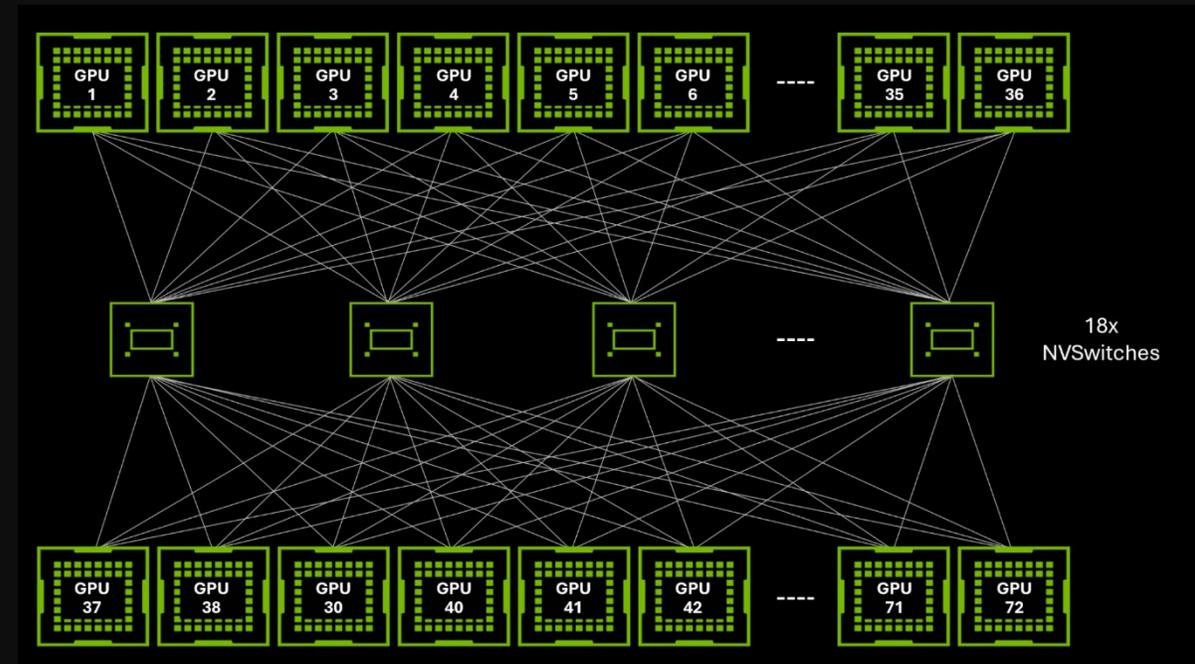
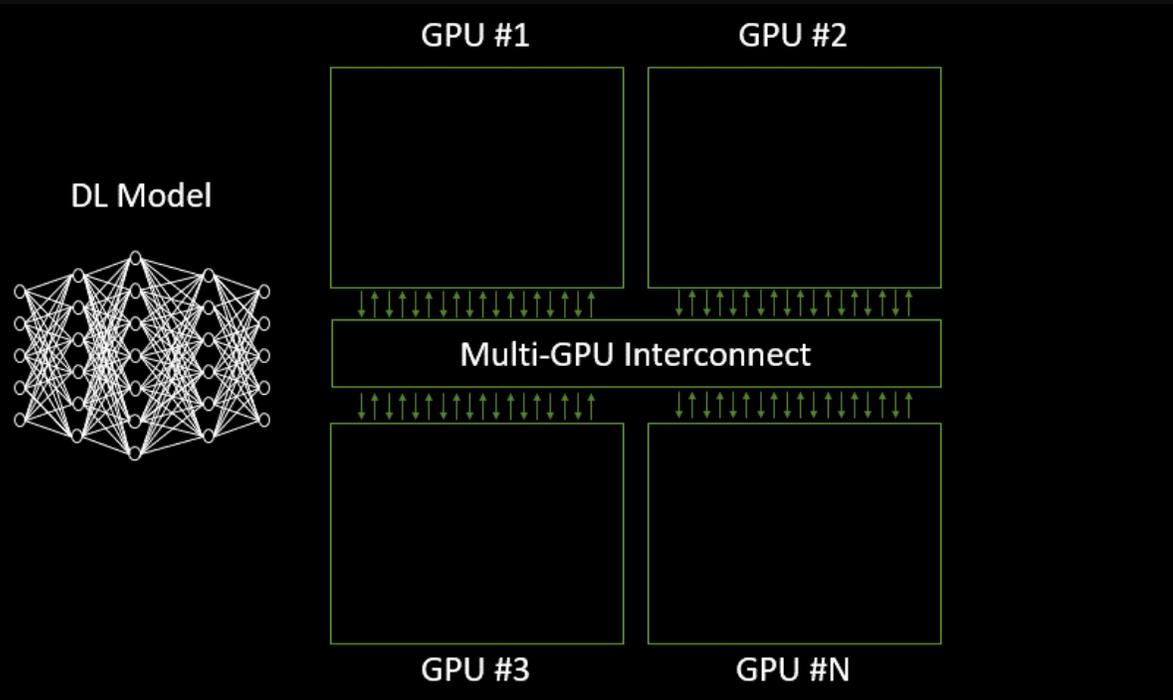
Rack-Scale Computing Platform



NVIDIA GB200 NVL72

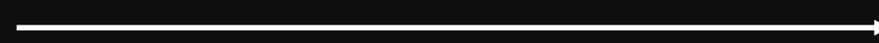
• ATTENTION

Multi-GPU Interconnect의 중요성



High volume of GPU-to-GPU communication happens with model parallelism

#N : 8



#N : 72

High-Bandwidth Domain

전력과 열관리의 필요성

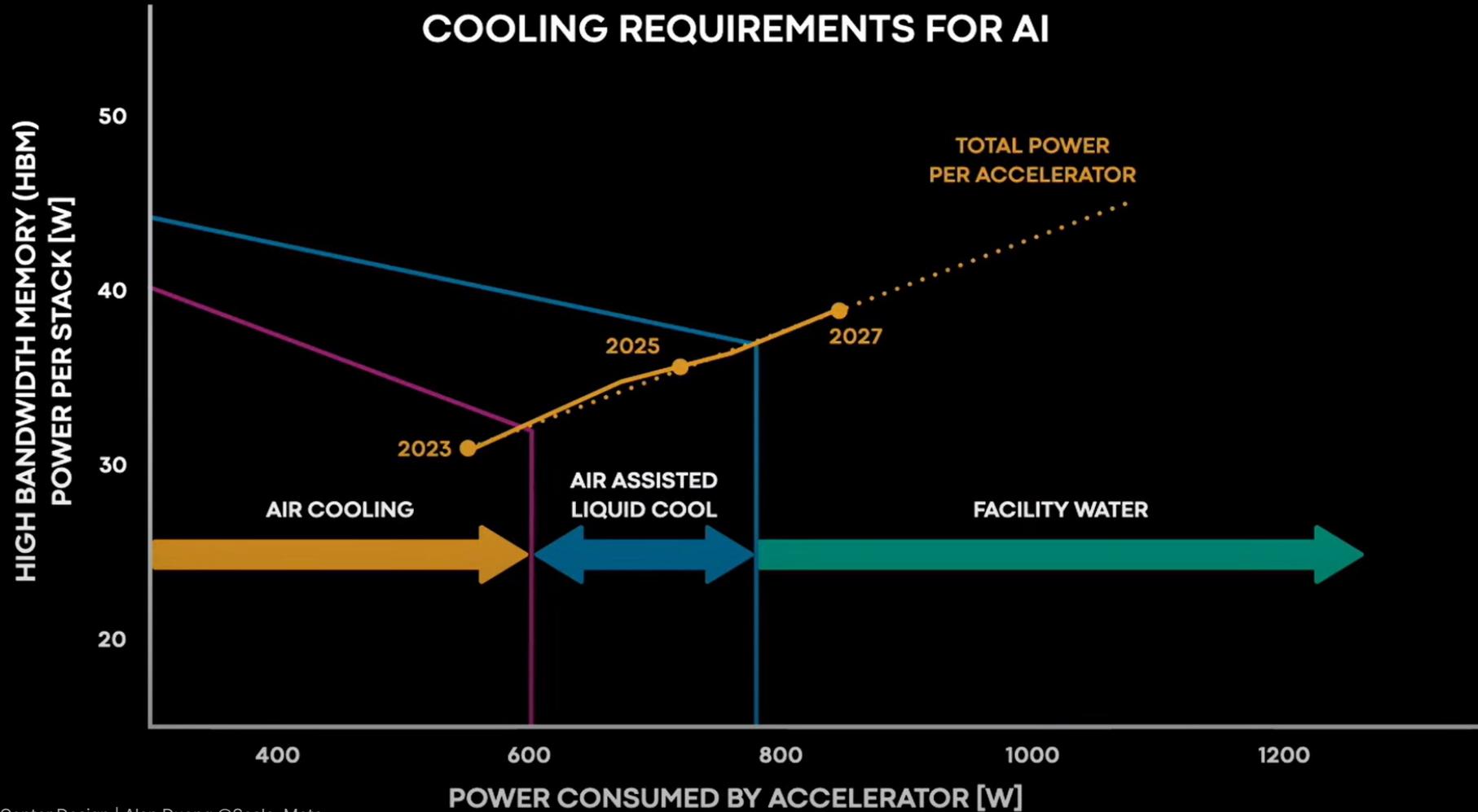
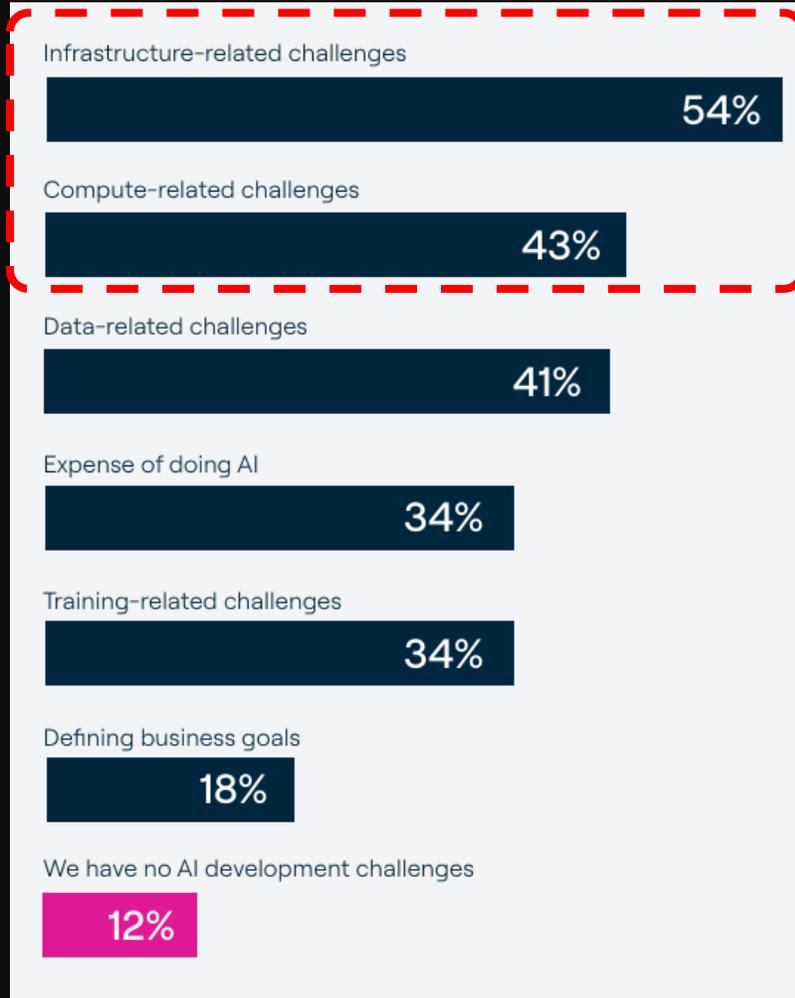


Image Source: Next-Generation Data Center Design | Alan Duong @Scale, Meta

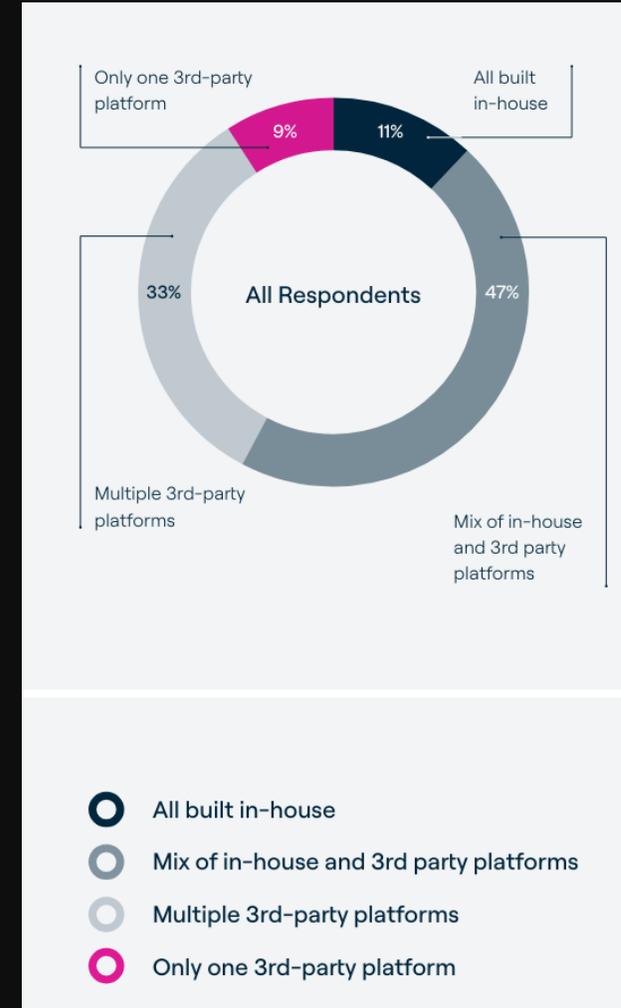
• ATTENTION

AI 인프라에 대한 어려움/복잡성

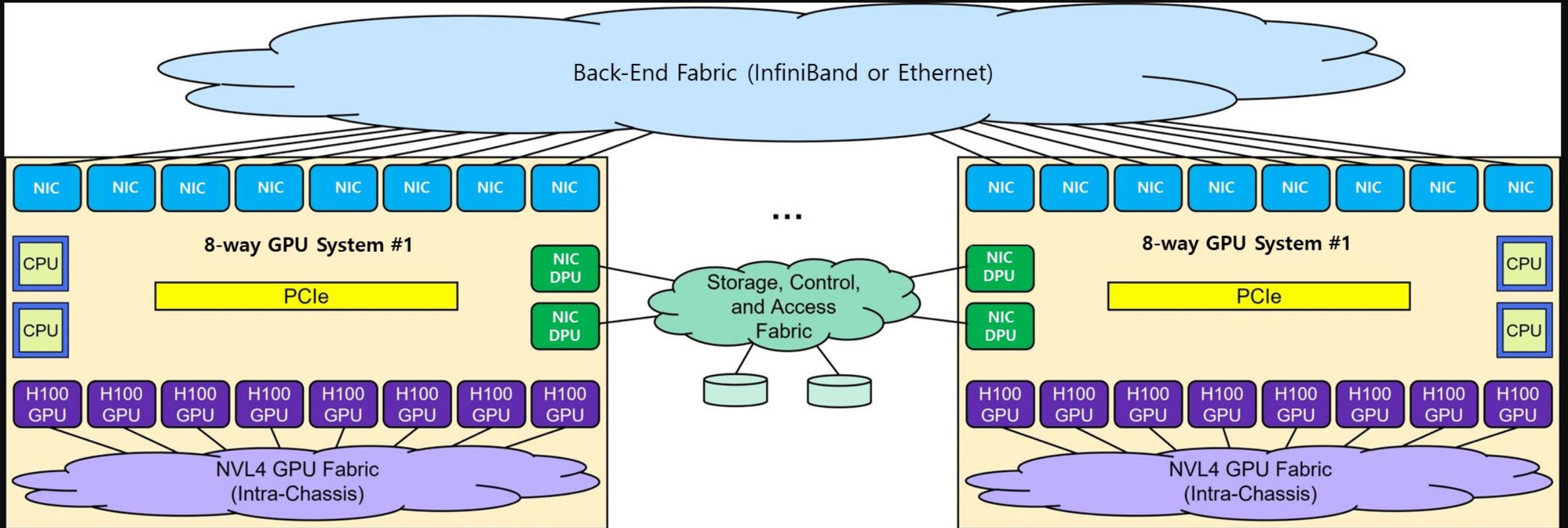
Challenges for AI Development



AI/ML Stack Architecture

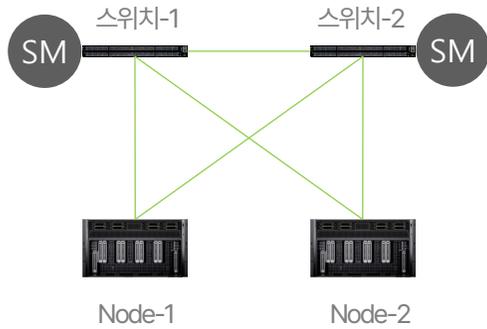


AI 인프라 (High Level) 네트워킹 모델

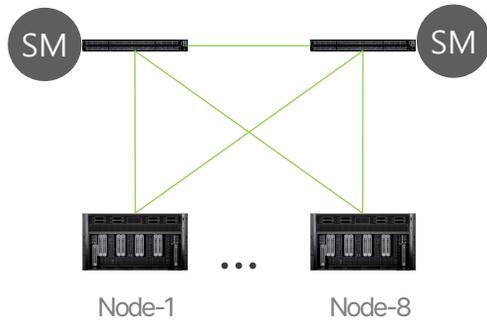


멀티노드 GPU 클러스터 기본 디자인

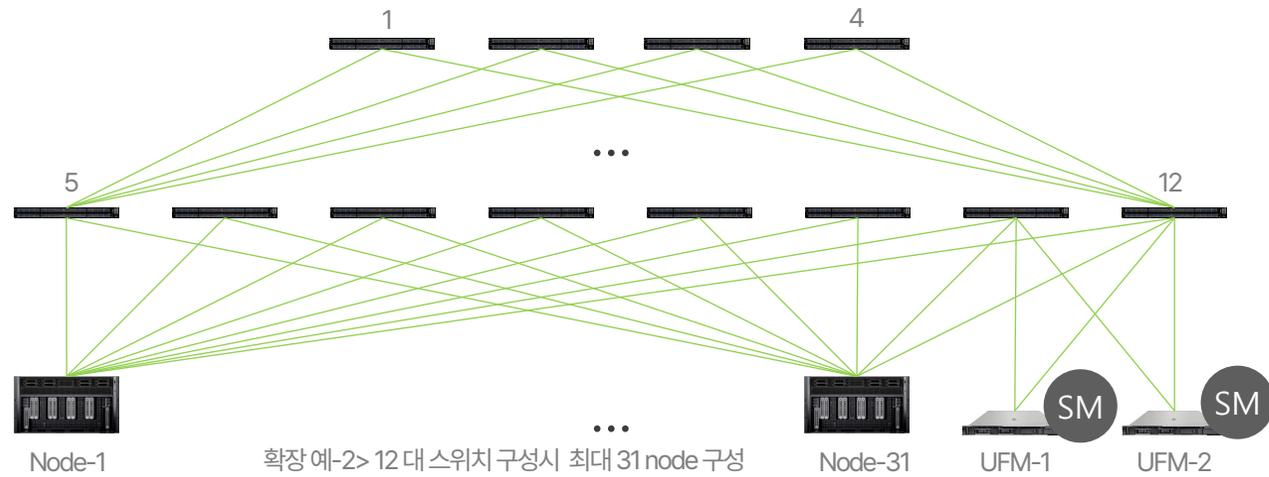
InfiniBand 확장 구성안



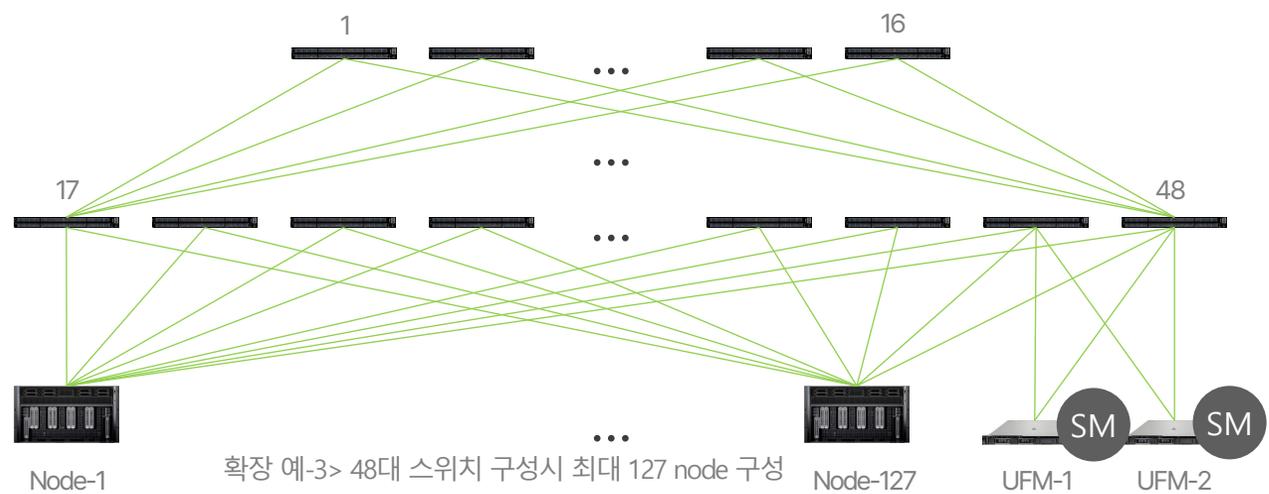
초기 2 node



확장 예-1> 최대 8 node 확장



확장 예-2> 12대 스위치 구성시 최대 31 node 구성



확장 예-3> 48대 스위치 구성시 최대 127 node 구성

- 서버넷 관리자(SM) – InfiniBand fabric 을 실행하고 관리하는 프로그램

• ATTENTION

GPU 클러스터 검증

새로운 GPU 클러스터를 구축한 경우, 시스템에 스트레스를 가하여 모든 장비가 올바르게 작동하고 동일한 성능을 발휘하는지 검증하는 단계를 반드시 거쳐야 합니다. DCGM 테스트, NCCL 테스트, HPL/MLPerf 벤치마크 테스트를 통해 클러스터의 각 노드가 동일한 수준의 계산/연산 성능 및 통신 성능을 수행하는지 검증합니다.



DCGM Test

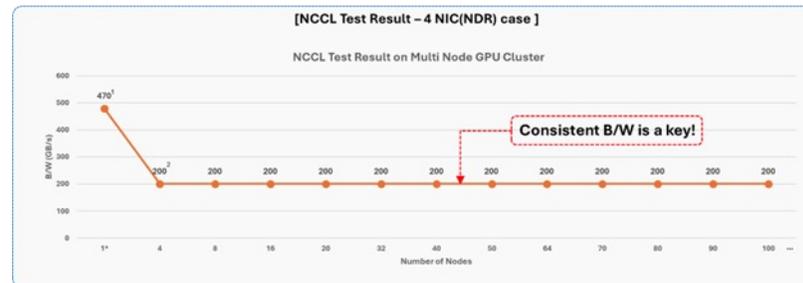
- 개별 Multi-GPU 시스템 기능 및 성능을 검증.

<pre>root@nccltest01:~# root@nccltest01:~# root@nccltest01:~# root@nccltest01:~# dcmi diag -r 4 Successfully ran diagnostic for group. ----- Diagnostic Result ----- Metadata DCGM Version 3.3.5 Driver Version Detected 545.23.08 GPU Device IDs Detected 2500,2506 Deployment DevList Pass NVML Library Pass CUDA Main Library Pass Partitions and OS Blocks Pass Persistence Mode Pass Environment Variables Pass Page Retirement/Row Ramp Pass Graphics Processes Pass Inform Pass ----- root@nccltest01:~#</pre>	<pre>root@nccltest01:~# root@nccltest01:~# root@nccltest01:~# dcmi nvlink --errors -g 0 ----- NVLINK Error Counts GPU 0 ----- root@nccltest01:~# dcmi nvlink --errors -g 1 ----- NVLINK Error Counts GPU 1 ----- root@nccltest01:~#</pre>
<p>dcmi diag -r 4 입력하여 하드웨어 진단 실시</p>	<p>dcmi nvlink --errors -g <gpu id> 입력하여 각 GPU의 nvlink 진단 실시</p>



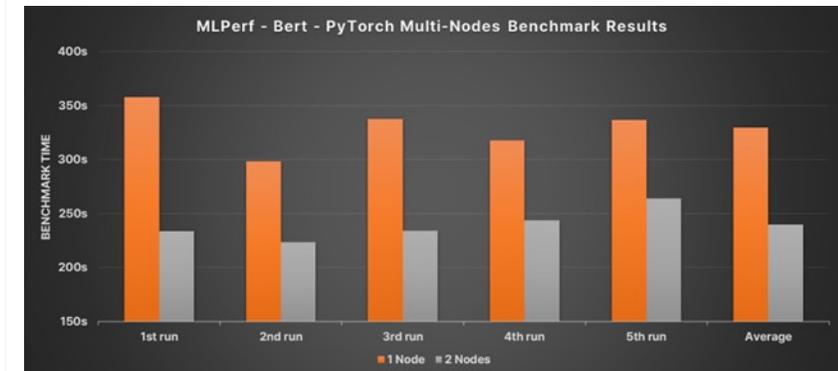
NCCL Test

- 노드 간 네트워크 대역폭이 예상대로 확보되는지 검증
- 노드 수가 증가함에 따라 네트워크 대역폭이 일정하게 유지되는 검증



HPL/MLPerf Test

- Multi-Node 환경에서 GPU Cluster의 계산/연산 성능이 예상대로 확보되는지 검증
- 노드 수가 증가함에 따라 계산/연산 성능이 일정하게 증가하는지 검증



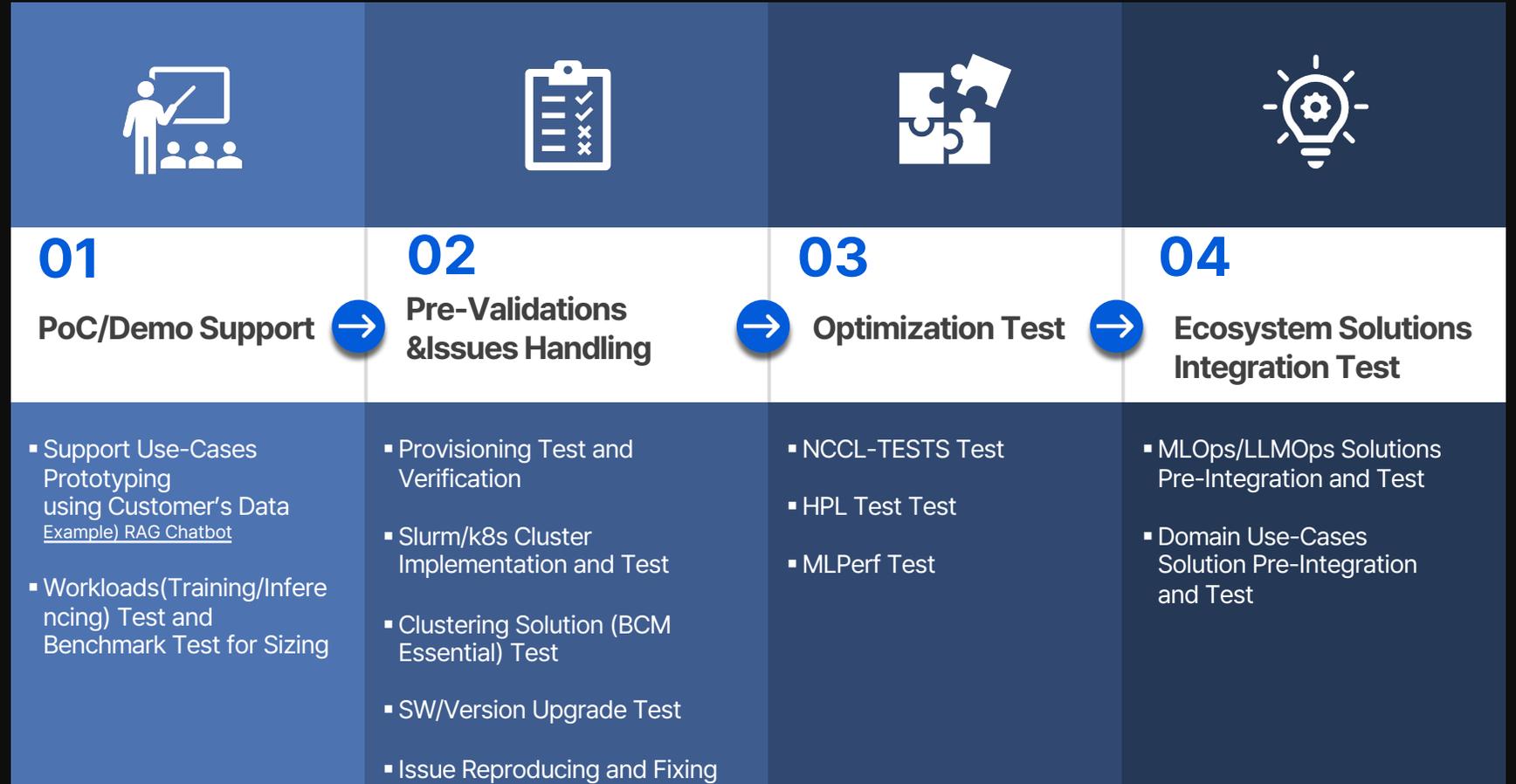
Dell AI Factory 플랫폼

메이머스트는 고객의 성공적인 AI 여정을 위해 Dell AI Factory 플랫폼을 운영하고 있습니다.



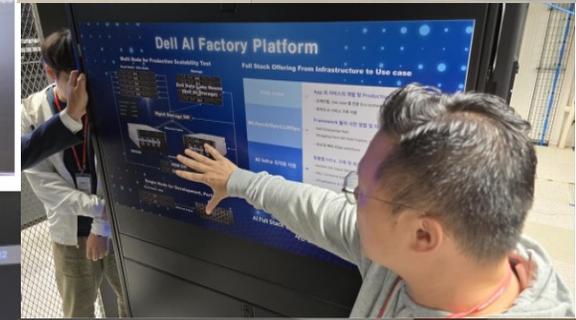
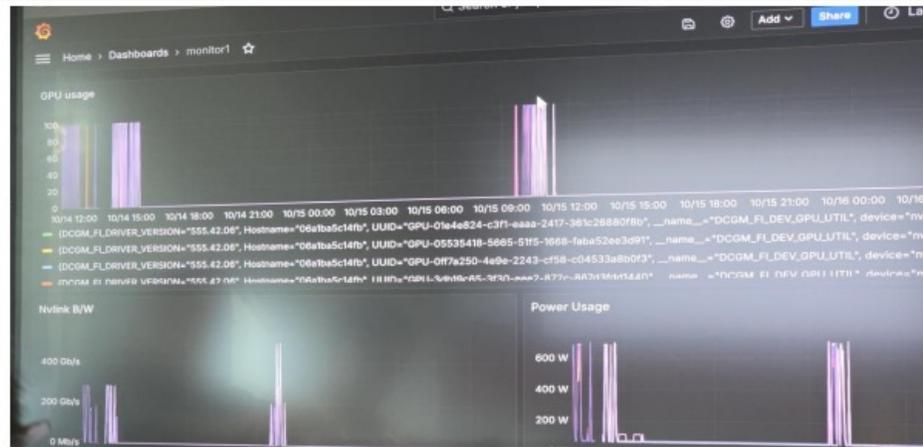
Daol TS | MUST

Dell AI Factory Demo/PoC Cetner



ATTENTION

AI Factory 플랫폼을 통한 고객 기술 지원, 시연 및 Benchmark



```

Device 0 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 6.395 GiB/s TX: 643.6 MiB/s
GPU 1860MHz MEM 2619MHz TEMP 73 FAN N/A POW 684 / 700 W 34.4236Gi/79.647Gi
GPU 1 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 3.645 GiB/s TX: 892.7 MiB/s
GPU 1086MHz MEM 2619MHz TEMP 62 FAN N/A POW 679 / 700 W 34.564Gi/79.647Gi
Device 2 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 6.207 GiB/s TX: 721.2 MiB/s
GPU 1935MHz MEM 2619MHz TEMP 68 FAN N/A POW 676 / 700 W 34.564Gi/79.647Gi
Device 3 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 2.122 GiB/s TX: 788.6 MiB/s
GPU 1808MHz MEM 2619MHz TEMP 72 FAN N/A POW 687 / 700 W 34.564Gi/79.647Gi
Device 4 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 5.666 GiB/s TX: 827.9 MiB/s
GPU 1845MHz MEM 2619MHz TEMP 74 FAN N/A POW 685 / 700 W 34.564Gi/79.647Gi
Device 5 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 1.922 GiB/s TX: 834.0 MiB/s
GPU 1875MHz MEM 2619MHz TEMP 62 FAN N/A POW 678 / 700 W 34.564Gi/79.647Gi
Device 6 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 6.525 GiB/s TX: 835.5 MiB/s
GPU 1830MHz MEM 2619MHz TEMP 70 FAN N/A POW 686 / 700 W 34.564Gi/79.647Gi
Device 7 [NVIDIA H100 80GB HBM3] PCIe GEN 5@16x RX: 2.428 GiB/s TX: 751.6 MiB/s
GPU 1860MHz MEM 2619MHz TEMP 68 FAN N/A POW 675 / 700 W 33.861Gi/79.647Gi
    
```

PID	USER	DEV	TYPE	GPU	GPU MEM	CPU	HOST MEM	Command
1015131	root	5	Compute	99%	34942MiB	43%	101%	3061MiB python -u /workspace/bert/run_pretraining.py --train_batch_size
1015150	root	2	Compute	99%	34942MiB	43%	101%	3043MiB python -u /workspace/bert/run_pretraining.py --train_batch_size
1015170	root	3	Compute	99%	34942MiB	43%	101%	3073MiB python -u /workspace/bert/run_pretraining.py --train_batch_size
1015154	root	4	Compute	99%	34942MiB	43%	101%	3073MiB python -u /workspace/bert/run_pretraining.py --train_batch_size
1015173	root	5	Compute	99%	34942MiB	43%	101%	3056MiB python -u /workspace/bert/run_pretraining.py --train_batch_size
1015155	root	6	Compute	99%	34942MiB	43%	101%	3045MiB python -u /workspace/bert/run_pretraining.py --train_batch_size

```

# Using devices
# Rank 0 Group 0 PID 1053880 on work001 device 0 [8x19] NVIDIA H100 80GB HBM3
# Rank 1 Group 0 PID 1053881 on work001 device 1 [8x20] NVIDIA H100 80GB HBM3
# Rank 2 Group 0 PID 1053882 on work001 device 2 [8x45] NVIDIA H100 80GB HBM3
# Rank 3 Group 0 PID 1053883 on work001 device 3 [8x5] NVIDIA H100 80GB HBM3
# Rank 4 Group 0 PID 1053884 on work001 device 4 [8x26] NVIDIA H100 80GB HBM3
# Rank 5 Group 0 PID 1053885 on work001 device 5 [8x26] NVIDIA H100 80GB HBM3
# Rank 6 Group 0 PID 1053886 on work001 device 6 [8x26] NVIDIA H100 80GB HBM3
# Rank 7 Group 0 PID 1053887 on work001 device 7 [8x26] NVIDIA H100 80GB HBM3
# Rank 8 Group 0 PID 1839542 on work002 device 1 [8x26] NVIDIA H100 80GB HBM3
# Rank 9 Group 0 PID 1839543 on work002 device 2 [8x45] NVIDIA H100 80GB HBM3
# Rank 10 Group 0 PID 1839544 on work002 device 3 [8x5] NVIDIA H100 80GB HBM3
# Rank 11 Group 0 PID 1839545 on work002 device 4 [8x26] NVIDIA H100 80GB HBM3
# Rank 12 Group 0 PID 1839546 on work002 device 5 [8x26] NVIDIA H100 80GB HBM3
# Rank 13 Group 0 PID 1839547 on work002 device 6 [8x26] NVIDIA H100 80GB HBM3
# Rank 14 Group 0 PID 1839548 on work002 device 7 [8x26] NVIDIA H100 80GB HBM3
# Rank 15 Group 0 PID 1839549 on work002 device 8 [8x26] NVIDIA H100 80GB HBM3
    
```

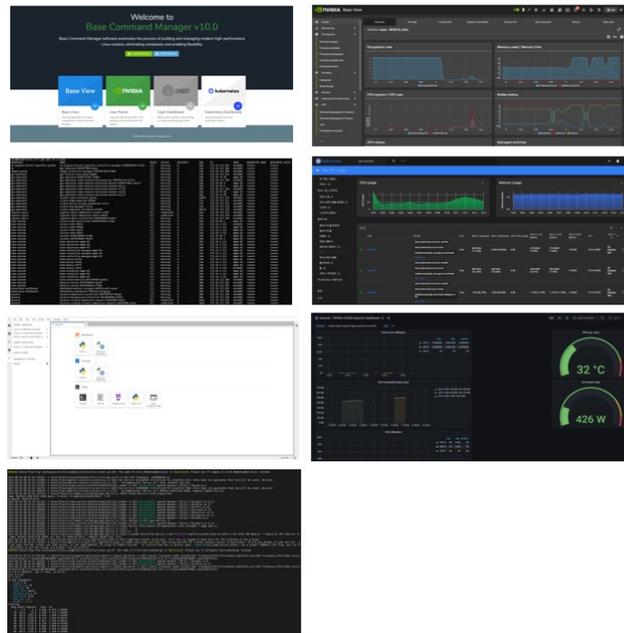
	size	count	type	redop	root	time (us)	algbw (GB/s)	busbw #wrong	time (us)	algbw (GB/s)	busbw #wrong	
(0)	(element)											
	1073741824	268225436	Float	sum	-1	4491.3	238.49	447.16	0	4491.7	239.85	448.22
	2147481840	538870912	Float	sum	-1	9016.3	238.18	446.58	0	9045.8	237.40	445.13
	4296987296	1073741824	Float	sum	-1	17232	240.25	407.34	0	17054	251.84	472.21
	8593974592	2147481824	Float	sum	-1	33331	237.72	403.22	0	33220	258.43	484.56



ATTENTION

AI Factory 플랫폼 활용 예시- 1/3

NVIDIA BCM10	
Test Plan	<ol style="list-style-type: none"> Environment configuration <ul style="list-style-type: none"> Head Node (Dell PowerEdge 760) BCM installation and configuration Provisioning Image Creation (Slurm, K8S) Master Node (Dell PowerEdge 760 3nodes) configuration and Worker Node (Dell PowerEdge XE9680 2nodes) Image Provisioning Slurm configuration and testing <ul style="list-style-type: none"> Slurm configuration through CLI commands Configuration verification and error checking GPU Stress Test K8S configuration and testing <ul style="list-style-type: none"> Configure K8S with CLI commands GPU Operator, Neteork Operator configuration Jupyter HUB, Grafana, Prometheus, K8S Dashboard, Ingress Control configuration Configuration verification and error checking GPU Stress Test
Result	<p>Nvidia bmc installation to ensure consistency of Dell infrastructure environments and consistency of key features related to k8 and GPUs</p> <ol style="list-style-type: none"> Slurm <ul style="list-style-type: none"> No special configuration or errors Verify normal operation of Multi GPU Stress K8S <ul style="list-style-type: none"> No special configuration or errors Verify normal operation of Multi GPU Stress <p>[Overall review] By configuring nvidia bcm, the customer can easily configure the necessary s/w stack, cluster the customer's gpu nodes, and monitor each node.</p>



TensorRT-LLM Inferencing	
Test Plan	<ol style="list-style-type: none"> Environmental Configuration <ul style="list-style-type: none"> OS: Ubuntu 22.04 BCM: 10.23.12 Docker: 24.09 Nvidia Container Toolkit: 1.16.0 Cuda Toolkit: 12.2 Model: Llama2 7B, Llama2 13B, Llama2 70B, Llama3 8B, Llama3 70B, Mistral 8*7B Planning <ul style="list-style-type: none"> Run TensorRT-LLM as Dockers Check GPU memory for each The same executable code applies different variables (8 GPUs, Batch Size, TP, InputLength, OutputLength). TFT, Latency, Through Measurement Reference URL <ul style="list-style-type: none"> https://github.com/NVIDIA/TensorRT-LLM/blob/v0.8.0/docs/source/performance.md
Result	<p>Each case by input length (Size from 1 to 1024)</p> <p>The batch size roughly doubles as the input length decreases by 1/2. As the input/output length decreases by 1/2, the batch size roughly doubles and the throughput roughly increases by 1.6-2.8x</p>

ISV Solutions	
Test Plan	<ol style="list-style-type: none"> Environment configuration <ul style="list-style-type: none"> OS: RHEL 8.7 K8S: 1.29 Master Node (Dell PowerEdge 760 3nodes), Work Node (Dell PowerEdge XE9680 2nodes) Plan <ul style="list-style-type: none"> 고객요구 사항 정의 사용자 계정, LB VIP, DNS, 스토리지 및 컨테이너 레지스트리 Ten AI PUB K8S Deploy Ten AI PUB Core Engine Deploy
Result	<p>Check out key features of K8, including consistency and GPU block splitting capabilities for Dell infrastructure environments.</p> <p>[Overall review] ai pub dev is composed of a UI so administrators can conveniently allocate resources. Multiple teams share resources and provide workload and monitoring capabilities. Its strength is that it provides advanced GPU management (GPU block splitting) functions and can split up to 1 GPU into 100.</p>

AI Factory 플랫폼 활용 예시- 2/3



<p>Test Plan</p>	<div style="text-align: right; background-color: #e91e63; color: white; padding: 5px; font-weight: bold;">Training /w NVIDIA Nemo FW</div> <p>1. Environment configuration</p> <ul style="list-style-type: none"> - OS: Ubuntu 22.04 - BCM: 10.23.12 - K8S: 1.27.13 - NeMo Framework Training container: 24.07 - GPU Operator: 24.6.2 - Network Operator: 23.5.0 - DataSet: databricks-dolly-15k - Llama 2 Model Download and convert model to NeMo <p>2. Plan</p> <ul style="list-style-type: none"> - Measurement of Fine-Tuning (SFT, Lora) time for Llama 2 models (7B, 13B, 70B) according to GPU quantity (4 GPU, 8 GPU) - Llama 2 7B, 4 GPU, Lora - Llama 2 7B, 8 GPU, (Lora, SFT) - Llama 2 13B, 4 GPU, Lora - Llama 2 13B, 8 GPU, (Lora, SFT) - Llama 2 70B, 8 GPU, Lora <p>3. Reference URL</p> <ul style="list-style-type: none"> - https://www.delltechnologies.com/asset/en-us/solutions/infrastructure-solutions/technical-support/gen-ai-model-customization-design-guide.pdf 																		
<p>Result</p>	<p style="text-align: center;">Comparison of LoRA and SFT Times for Different Models and GPU Configurations</p> <table border="1"> <thead> <tr> <th>Model (Number of GPUs)</th> <th>LoRA (time in minutes)</th> <th>SFT (time in minutes)</th> </tr> </thead> <tbody> <tr> <td>Llama 2 7B (4 x H100 SXM)</td> <td>81</td> <td>-</td> </tr> <tr> <td>Llama 2 7B (8 x H100 SXM)</td> <td>45</td> <td>107</td> </tr> <tr> <td>Llama 2 13B (4 x H100 SXM)</td> <td>219</td> <td>-</td> </tr> <tr> <td>Llama 2 13B (8 x H100 SXM)</td> <td>113</td> <td>230</td> </tr> <tr> <td>Llama 2 70B (8 x H100 SXM)</td> <td>892</td> <td>-</td> </tr> </tbody> </table> <p>[Overall review] Fine tuning (Lora, SFT) time was measured for Llama 2 7B, 13B, and 70B models with XE9680 H100 8 GPU. Fine tuning from 7B to 70B (excluding SFT) is possible with XE9680 1 Node.</p>	Model (Number of GPUs)	LoRA (time in minutes)	SFT (time in minutes)	Llama 2 7B (4 x H100 SXM)	81	-	Llama 2 7B (8 x H100 SXM)	45	107	Llama 2 13B (4 x H100 SXM)	219	-	Llama 2 13B (8 x H100 SXM)	113	230	Llama 2 70B (8 x H100 SXM)	892	-
Model (Number of GPUs)	LoRA (time in minutes)	SFT (time in minutes)																	
Llama 2 7B (4 x H100 SXM)	81	-																	
Llama 2 7B (8 x H100 SXM)	45	107																	
Llama 2 13B (4 x H100 SXM)	219	-																	
Llama 2 13B (8 x H100 SXM)	113	230																	
Llama 2 70B (8 x H100 SXM)	892	-																	

<p>Test Plan</p>	<div style="text-align: right; background-color: #e91e63; color: white; padding: 5px; font-weight: bold;">RAG prototyping /w Customer's data</div> <p>1. Environment configuration</p> <ul style="list-style-type: none"> - OS: Ubuntu 22.04 - BCM: 10.23.12 - K8S: 1.27.13 - GPU Operator: 24.6.2 - Network Operator: 23.5.0 - DataSet: XE9860, XE8640 technical-guide - Pre-trained Model: Mistral 7B instruct <p>2. Plan</p> <ul style="list-style-type: none"> - K8s Rag Chatbot Container configuration - Data set creation and labeling in XE9860, XE8640 technical-guide - Fine tuning (check accuracy) with Mistral 7B model data set - Writing Huggingface Pipeline Code - Code customization <p>3. Reference URL</p> <ul style="list-style-type: none"> - https://github.com/dell-examples/generative-ai
<p>Result</p>	<p>When asking the chatbot about servers on xe9680, 8 out of 10 questions were answered correctly.</p> <p>[Overall review] After selecting a base model with good performance, creating a dataset of the customer's domain area and the customer's in-house data, and then fine tuning the model, you can configure a RAG chatbot specialized for the customer.</p>

AI Factory 플랫폼 활용 예시 - 3/3

MayMust x + Frontier-Class Models Test

127.0.0.1:7860

MAYMUST Retrieval Digital Assistant - Multimodal

Model: nvidia/nvLM-D-72B

Load Model

GPU	mem(mib)
0	35000
1	25000
2	25000
3	25000
4	25000
5	25000
6	25000
7	25000

GPU	Usage(%)
0	10
1	10
2	10
3	10
4	10
5	10
6	10
7	10

Chat

Chatbot

who is the person in the image?

The person in the image is Dr. Jen-Hsun Huang, a Taiwanese-American businessman and engineer. He is the co-founder, president, and CEO of NVIDIA, a leading company in graphics processing units (GPUs) and artificial intelligence technology. Dr. Huang has played a significant role in advancing computer graphics and has been recognized for his contributions to the technology industry.

Powered by DELL PowerEdge XE9680 NVIDIA NVIDIA H100

AI Full Stack을 위한 에코시스템

AI Enterprise Ecosystem Solutions

TEN

crowdworks

upstage

XIIIab [씨이랩]

MINDs
company

MakinaRocks

TESTWORKS

AiNEWT

nVIDIA

labup

42MARU

Human & AI
RAINBOW BRAIN

Red Hat

SK telecom

Optimized & Validated GenAI Full-Stack Infra.

MAYMUST

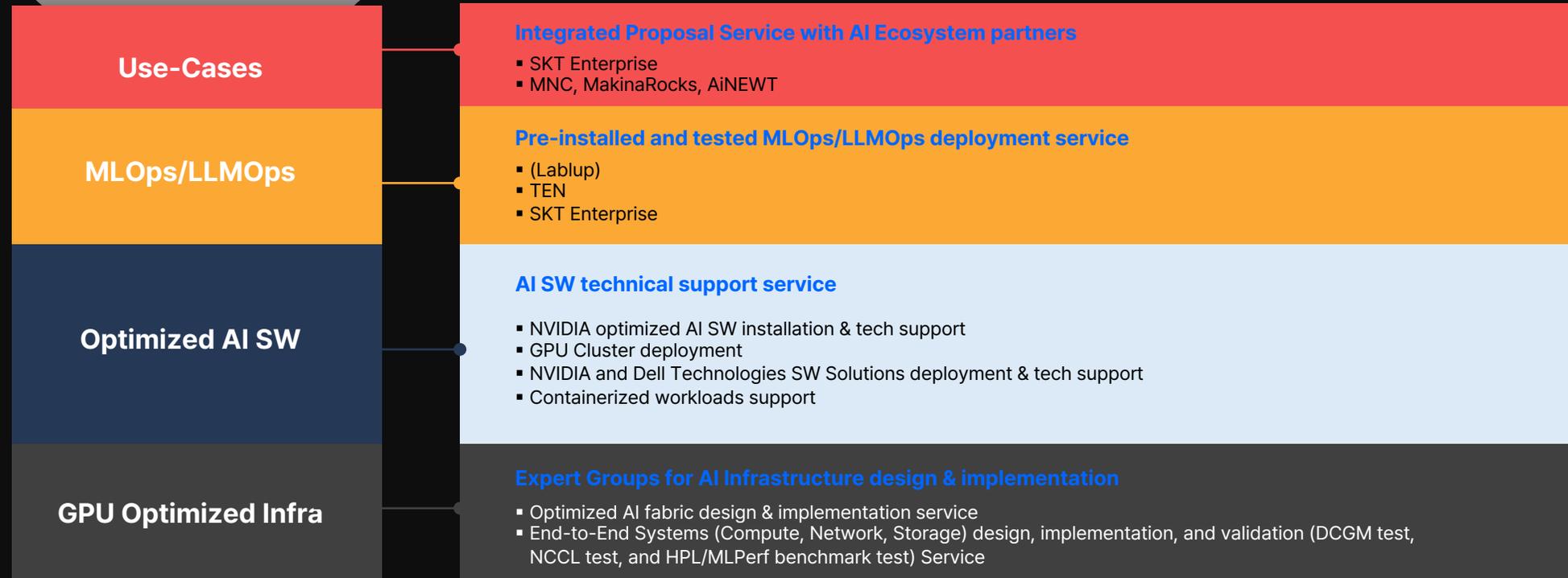
고객의 빠른 성공을 위해, AI Full Stack을 제공해 드립니다.

메이머스트의 최우선 목표는 고객의 성공입니다.

업계 최고의 에코시스템 파트너사와 함께 고객의 다양한 AI 워크로드 및 Use-Case를 최적의 인프라 환경하에서 지원합니다.

Manufacturing	Government	Healthcare
Education	Military	Entertainment
Finance	Retail	...

Full Stack Offering from Infra, SW to Use-case



THANK **YOU!**