



# Advancing Memory and Storage Architectures for Next-Gen AI Workloads

Vikram Sharma Mailthody

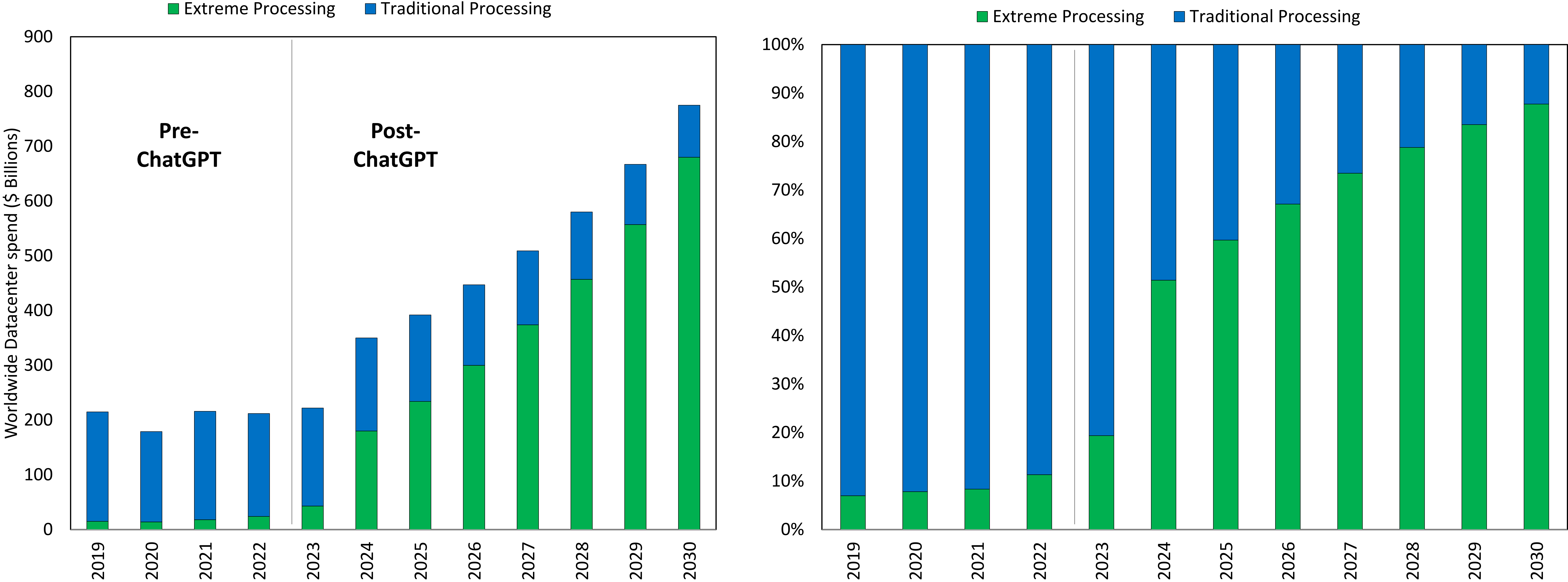
Sr. Research Scientist, NVIDIA Research





# Datacenter Compute Inversion

Extreme processing using GPUs and accelerators is now a norm in datacenter



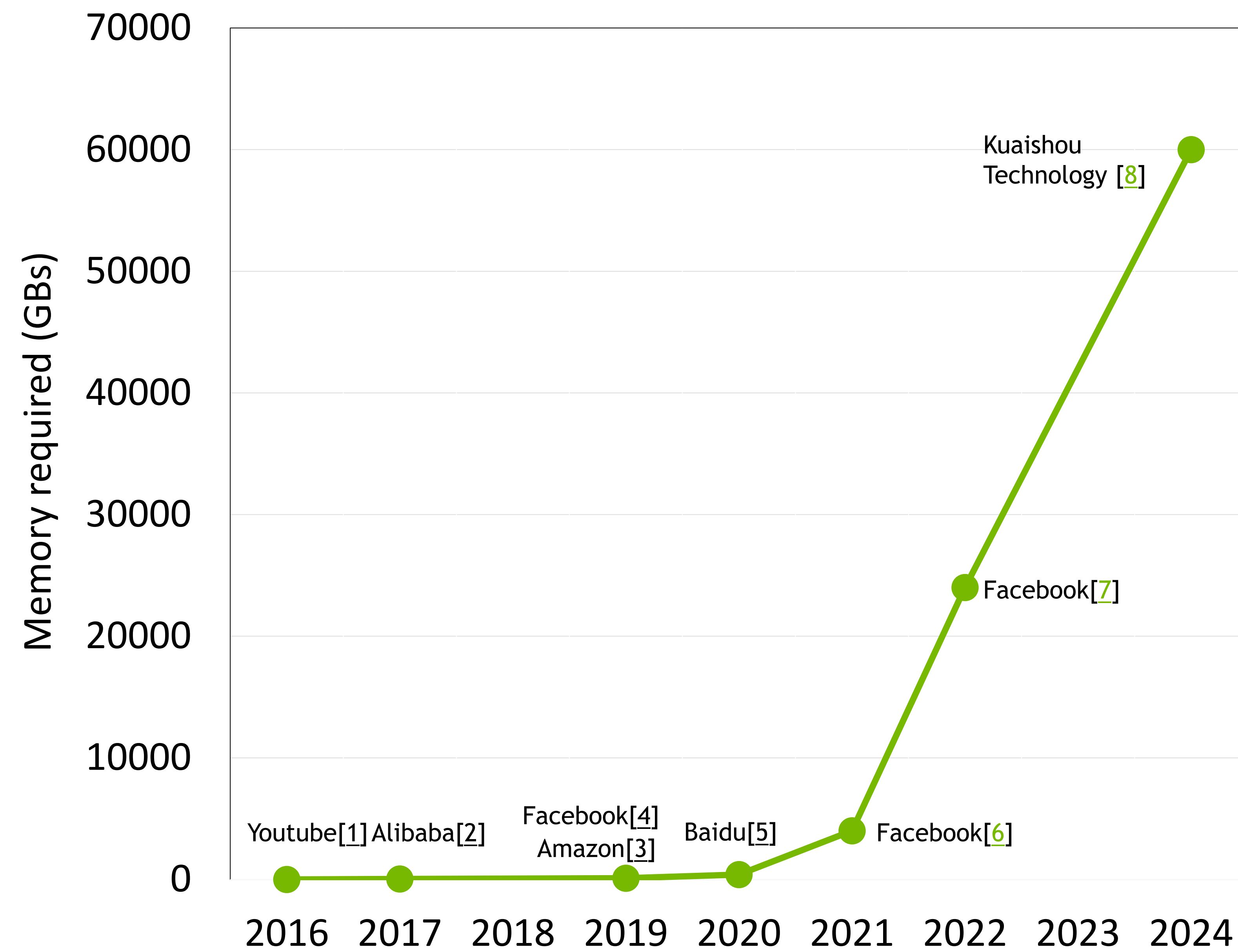
Data source: theCUBE Research 2025, Research & Visualization by David Floyer

# Emerging Generation And Prediction Based Applications

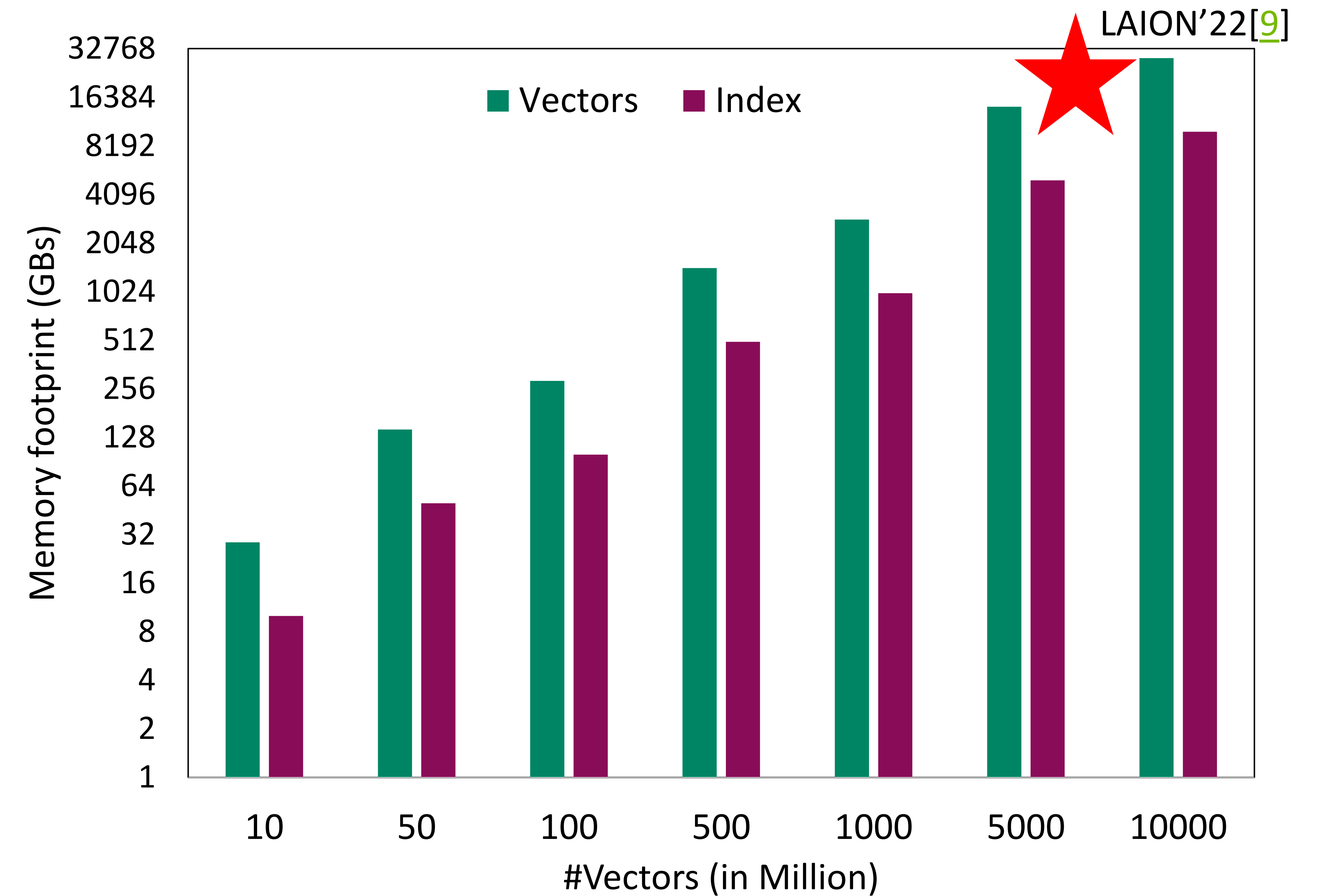
Prediction opens up new venues to revenue

	Generation	Prediction
Output	Query response Generated code Art, e.g. designs, music	Recommendation Filtered selection TopK from vector search
Usages	Chatbots, reasoning	eCommerce, social media, enterprise market
Cost of error	Wrong, hallucinated answer Debugging, escaped bugs Artistic artifact	Mispredictions ok if usually right
Infrastructure	LLM	Two tower models, GNN, relational graph transformer, Prefill heavy LLM inference
Computational complexity	High	Modest
Data needed in memory	Model weights	Knowledge/relational graph/embeddings
Data needed in storage	KV\$	Same, for bigger problems
Challenges	Running out of public data	Retrieval heavy, Many new horizons ("big data")

# Emerging workloads demand memory



Recommender Systems  
1TB to 100TBs



Vector databases  
1TB to 100TBs

# AI application overview

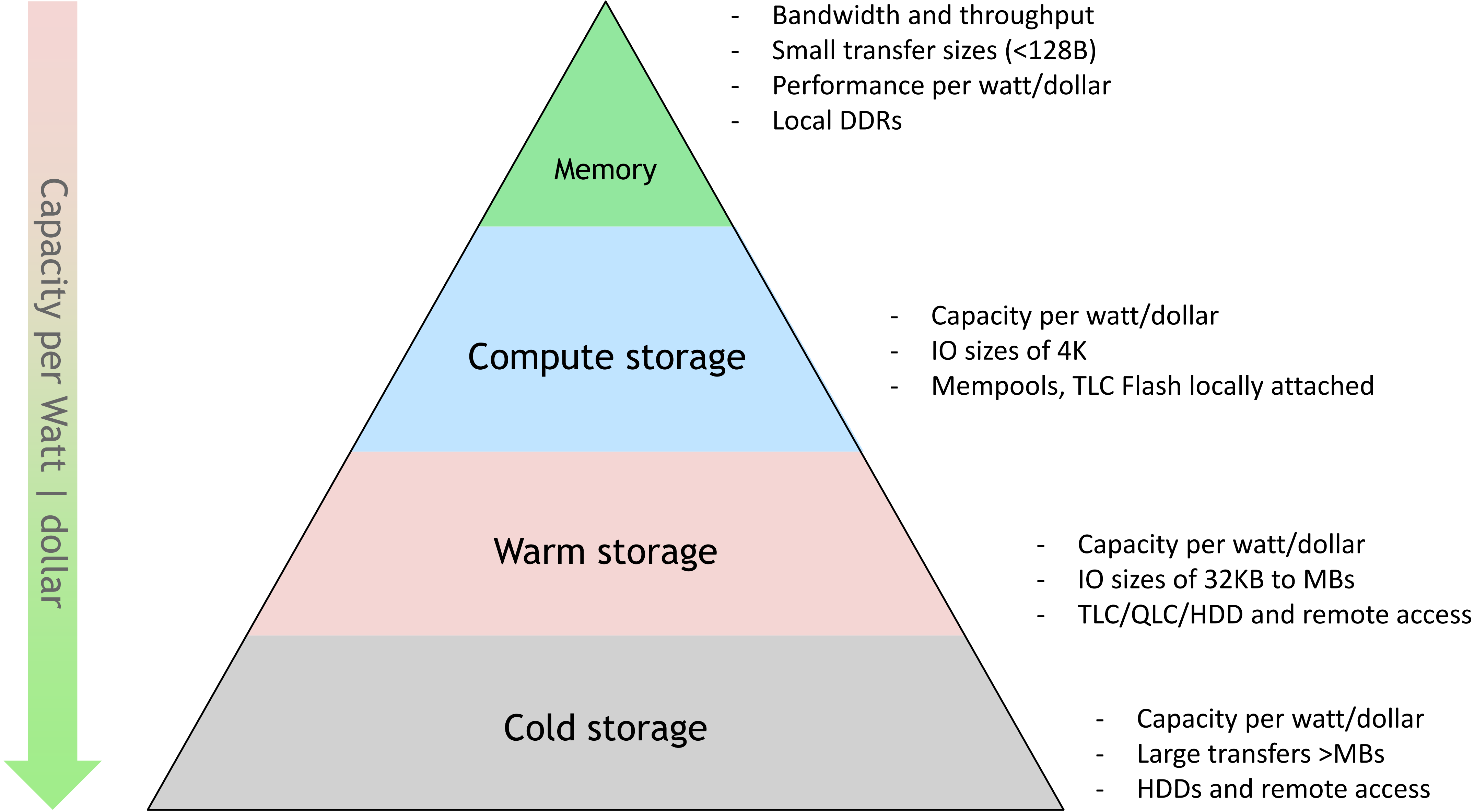
Apps bifurcate by access pattern and IO intensity; TB/TCO persists, IOPS/TCO is emerging

Area	Usage model	Applications	Access granularity	Total size /worker
Training	Checkpoint save/restore	LLM pretraining, fine tuning	10MB – 1sGB	1-10TB
Inference	KV context caching across queries, docs	LLM inference	8KB – 4MB	>10sTB
	LLM+GNN, GNN+LLM	Contextual LLMs	512B – 8KB	5TB – 400TB
	Vector database	Dynamic Index build	64B – 4KB	6.4Gb – 20TB
		LLM RAG doc retrieval	512B – 8KB	400GB – 1PB
Predictive AI		Graph RAG	64B – 8KB	400GB – 1PB
	Recommenders	64B – 4KB	5TB – 400TB	
	GNN induced subgraphs	eCommerce, fraud, social networks	512B – 8KB	>2TB
	Anomaly detection	eCommerce, fraud, social networks	512B – 8KB	>10TB
	Relational graphs	Data Science Automation	8B – 4KB	>100sTBs



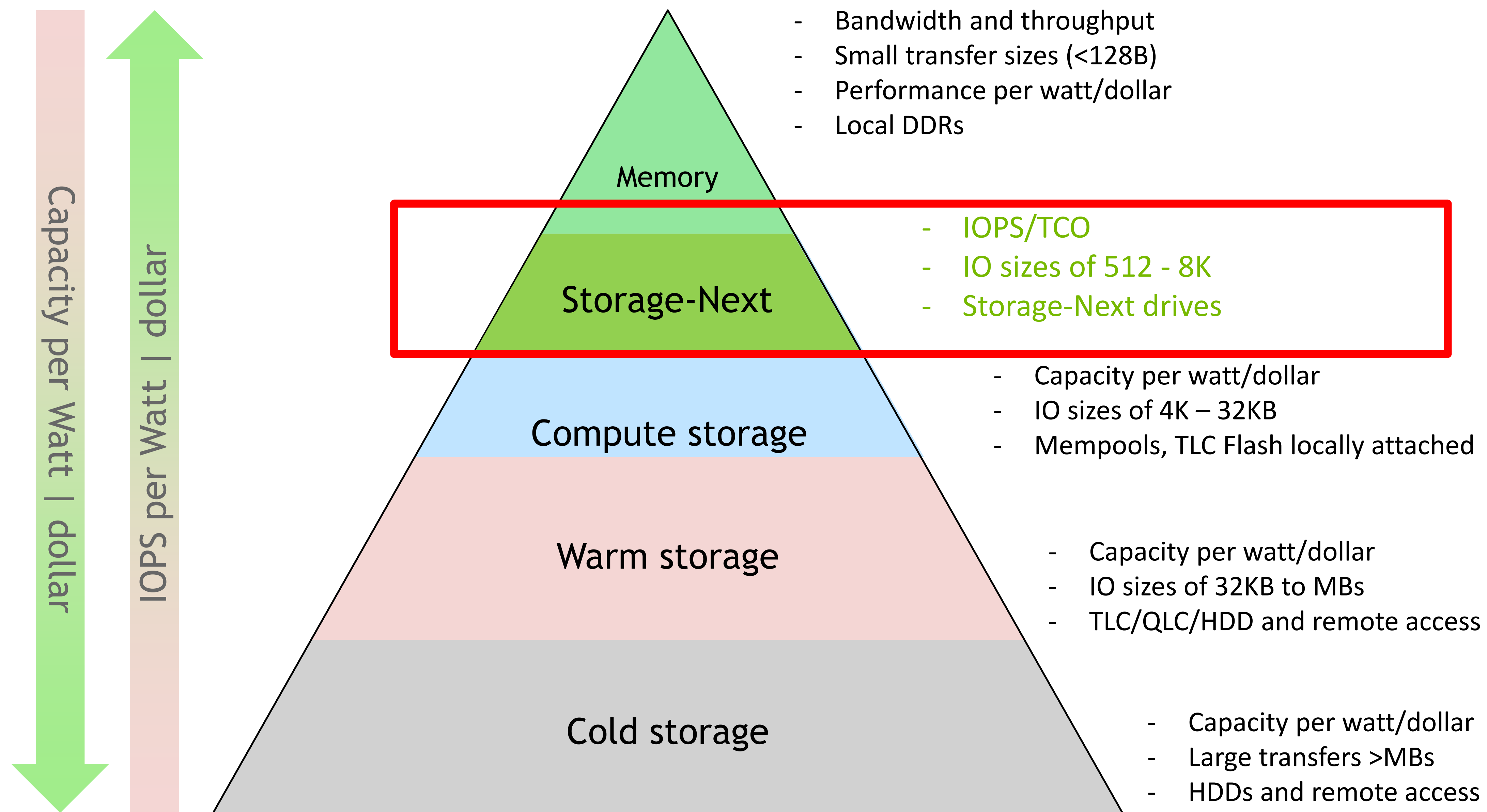
# Data hierarchy

Pre-GenAI era



# Data hierarchy

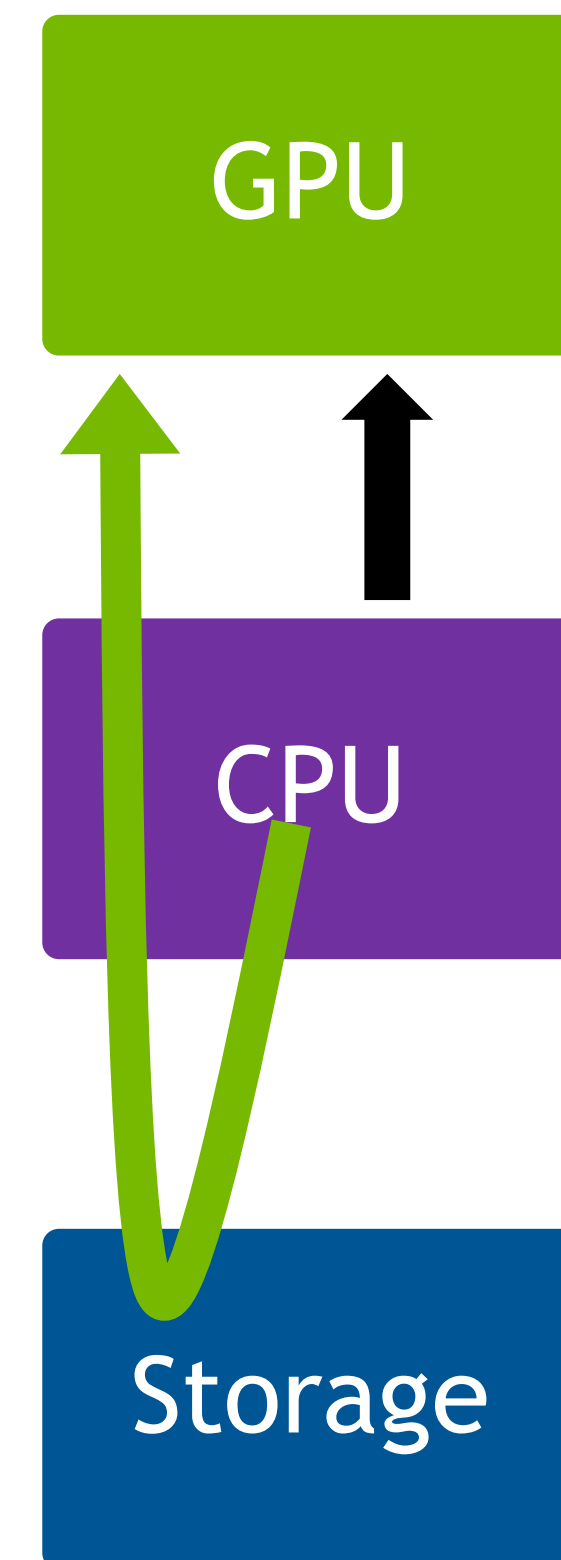
## Apparent Post-GenAI era



Trends existed in early 2016-2017

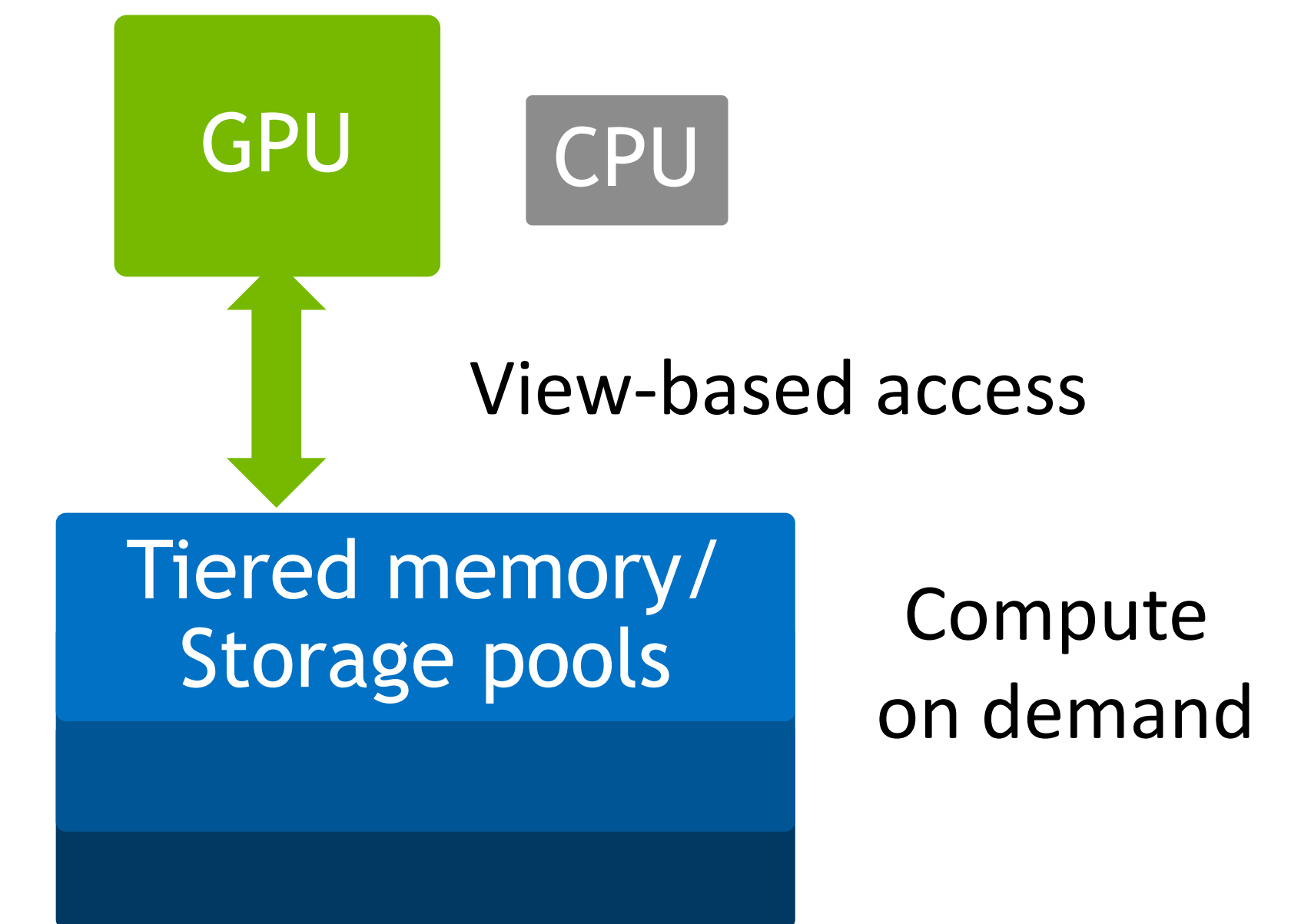
# From Offload Devices to Orchestrators

## Rethinking the Accelerator-Data Interface



Current Approach

- Current Approach
  - CPU has entire control path
  - GPU used as assistant
  - But majority of the work happens in GPU
  - Inefficient to load PB of data via tiling
- What applications need
  - A view-based access with a notion of infinite memory pool
  - Ability to fetch only needed data on-demand during computation
  - Ability to scale compute and data pipelines independently
  - Move majority of control to GPU while using CPU for house keeping



What applications need

Can GPU sustain enough parallelism to hide storage memory latency and provide a tiered memory/storage pool for applications?



# Software and Storage are the new bottleneck!

Little's Law

$$Q_d = T * L$$

Minimum steady state queue depth = Throughput \* Latency

- PCIe x16 Gen6 = 104GBps
  - For 512B access :—  $T = 104\text{GBps}/512\text{B} = 208 \text{ M IOPS}$
  - For 4KB access :—  $T = 104\text{GBps}/4\text{KB} = 26\text{M IOPS}$
- Assume SSD average access latency = 100us
  - $Q_d$  for 512B =  $208 \text{ M} * 100\text{us} = 20,800$
  - $Q_d$  for 4KB =  $26 \text{ M} * 100\text{us} = 2600$

GPUs and emerging workloads have enough parallelism to issue these many requests in-flight<sup>1</sup>.

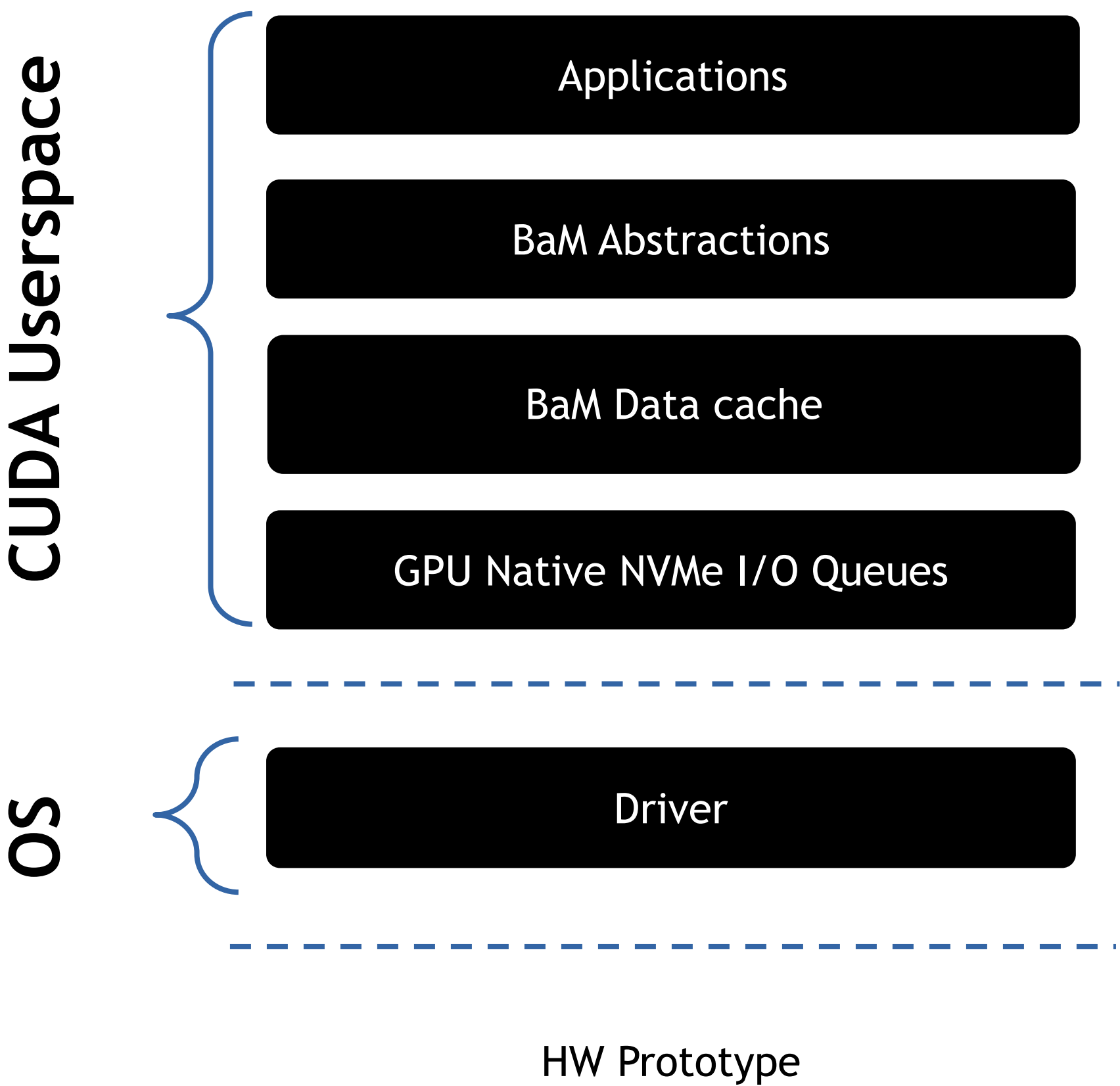
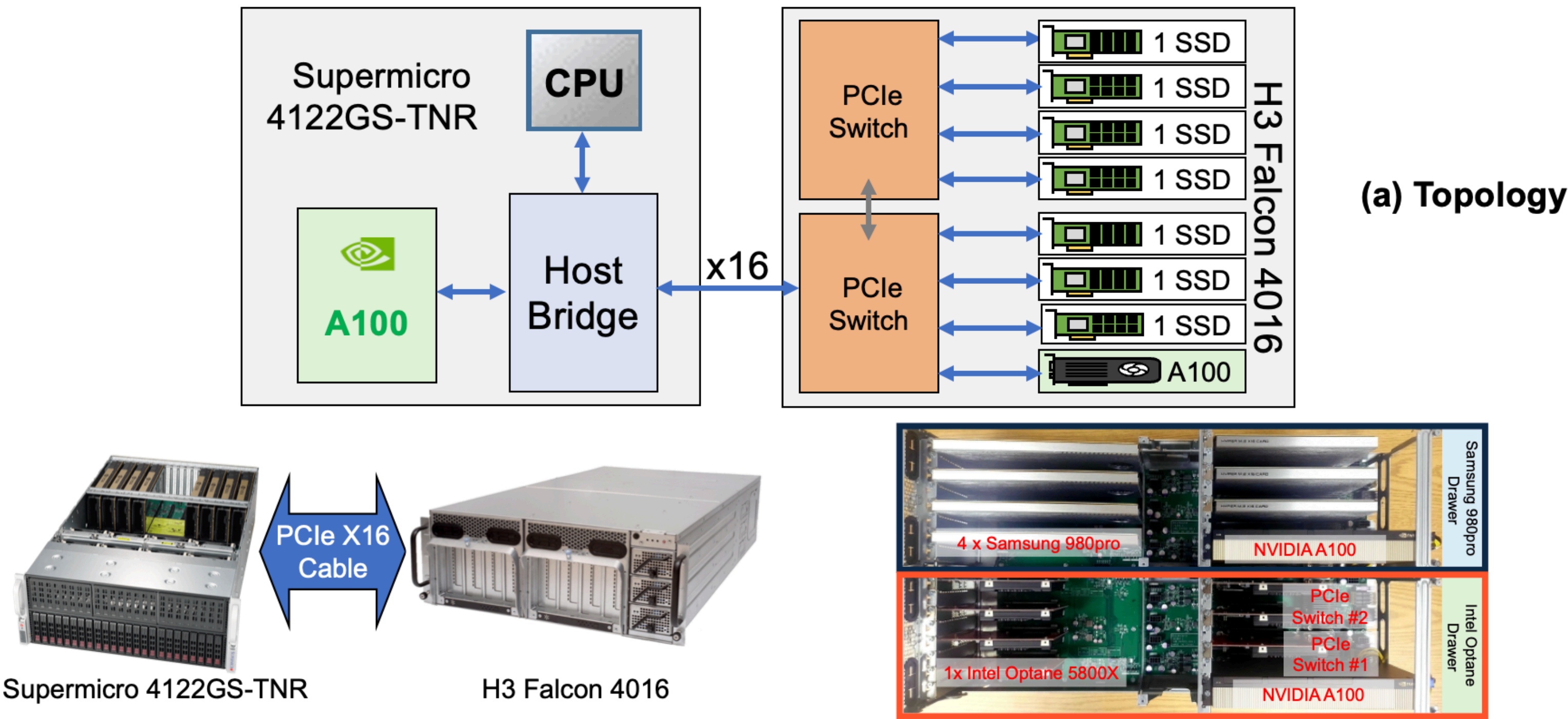
But the software stack and SSDs can't keep up.

CPU-drive software serialize, batch, or block effectively limiting QD resulting in underutilized bandwidth and stalling accelerators



# BaM

After many failures we got proof of concept prototype



**Large Graphs**

1.0X for BFS  
1.5X for CC

4.6X Perf/\$ for BFS  
6.6X Perf/\$ for CC

Name	Team	Number	Position	Age	Height	Weight	College	
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas
1	Joe Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State
4	Jonas Jerebko	Boston Celtics	8.0	PF	28.0	6-10	231.0	N/A
5	Amir Johnson	Boston Celtics	90.0	PF	29.0	6-9	240.0	N/A
6	Jordan Mickey	Boston Celtics	55.0	PF	21.0	6-8	239.0	LDSU
7	Kelly Olynyk	Boston Celtics	41.0	C	25.0	7-0	238.0	Georgia
8	Terry Rucker	Boston Celtics	12.0	PG	22.0	6-2	190.0	Louisville
9	Marcus Smart	Boston Celtics	35.0	PG	22.0	6-4	220.0	Oklahoma State
10	Jared Sullivan	Boston Celtics	7.0	C	24.0	6-9	260.0	Ohio State
11	Isiah Thomas	Boston Celtics	4.0	PG	27.0	5-9	185.0	Washington
12	Evan Turner	Boston Celtics	11.0	SG	27.0	6-7	220.0	Ohio State
13	James Young	Boston Celtics	13.0	SG	20.0	6-6	215.0	Kentucky
14	Tyler Zeller	Boston Celtics	44.0	C	26.0	7-0	253.0	North Carolina

**Data Analytics**

Upto 5.3X perf

Upto 23.6X Perf/\$

**GNN**

Upto 8.3X perf

Upto 38.3X Perf/\$



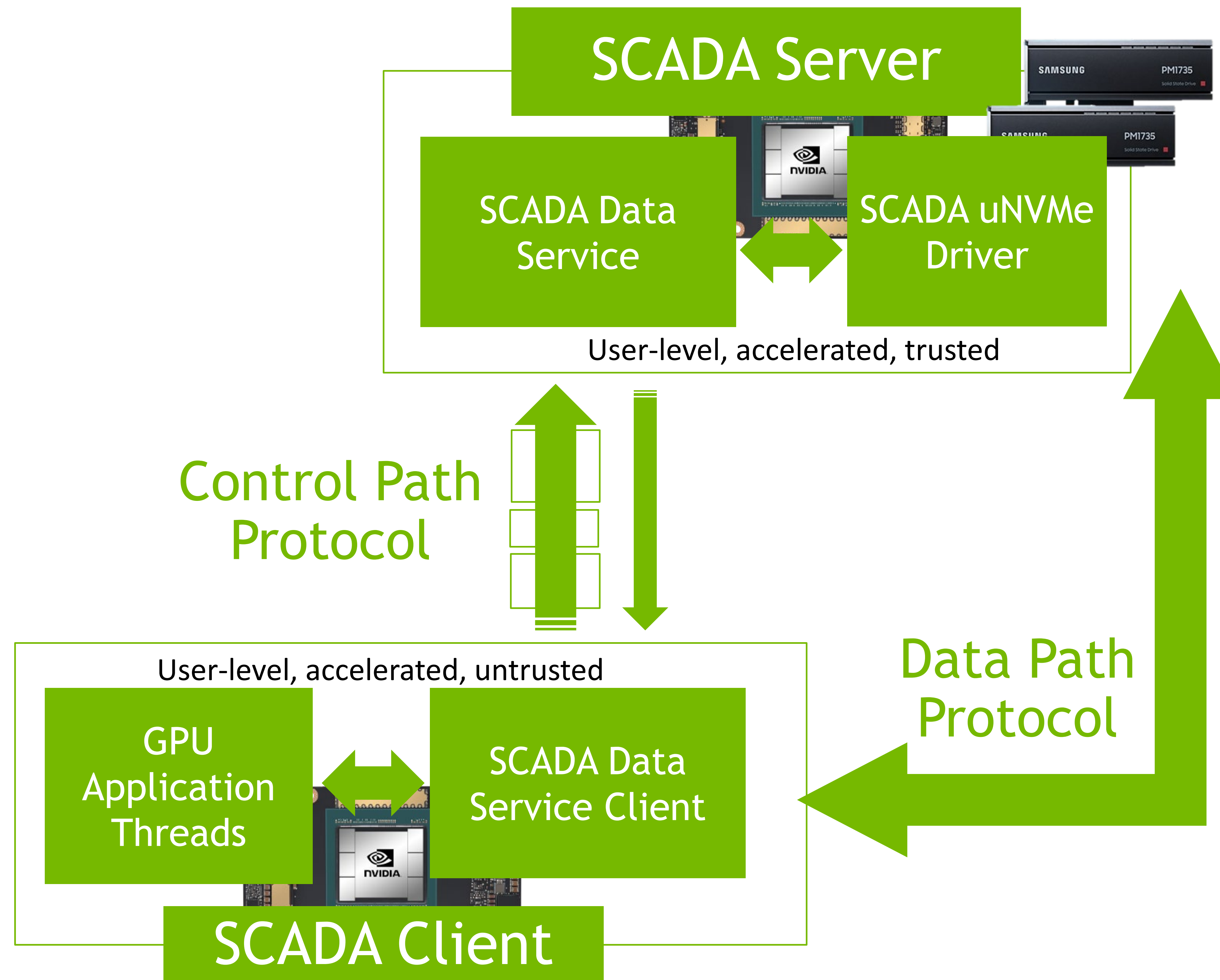
Many more...

But ...



# SCADA – SCaled Accelerated Data Access

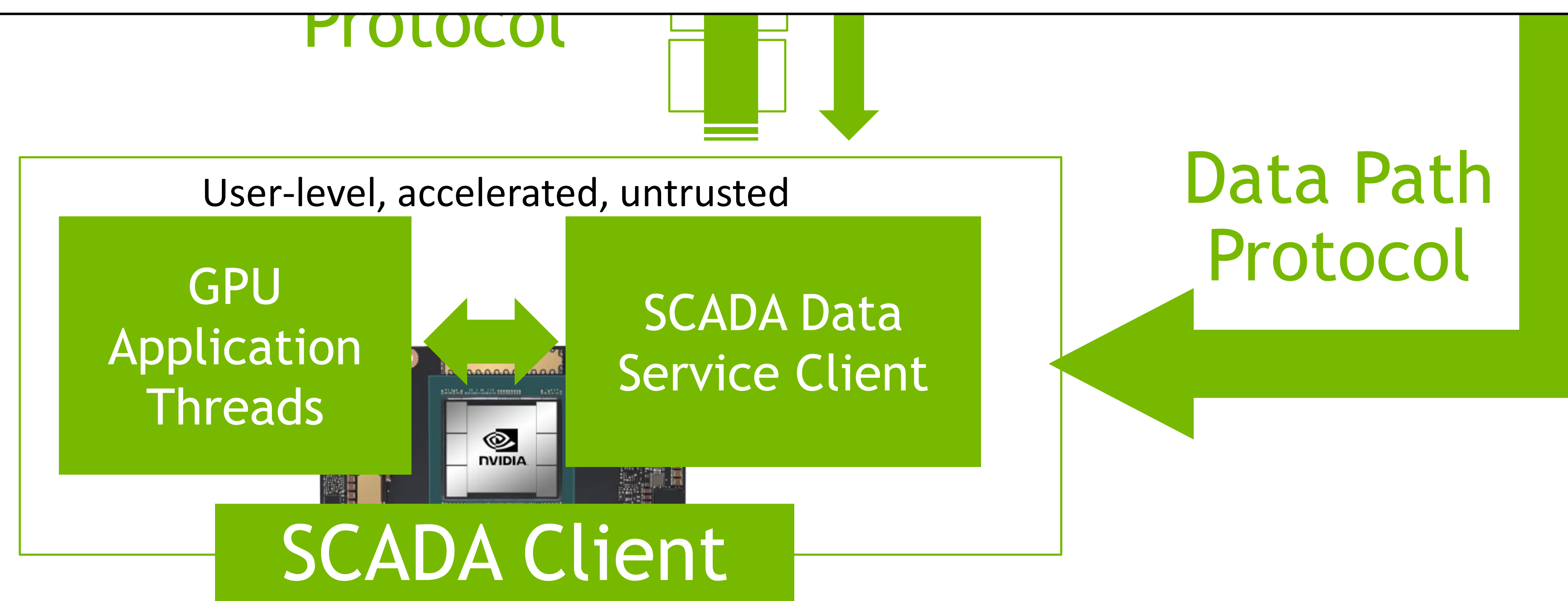
Software Architecture and Components – culmination of 8+ years of research and engineering



# SCADA – Scaled Accelerated Data Access

## Software Architecture and Components

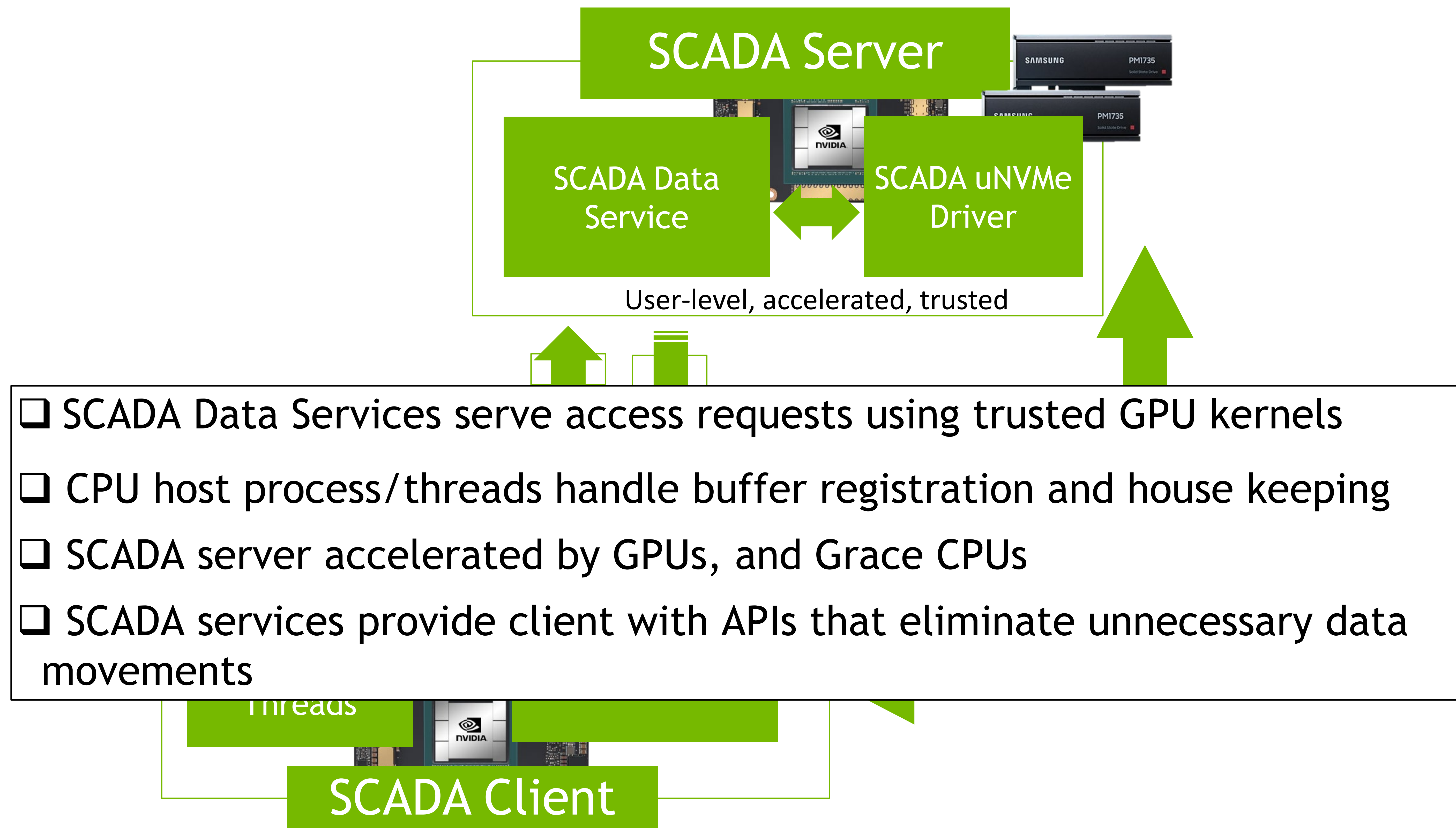
- ❑ Data Service Client is a header-only library compiled with applications
- ❑ Accelerate data buffer management through application defined (customized) software cache in HBM
- ❑ Provides simple abstractions that works with thrust library
  - ❑ mdspan array, keyvalue, graph, ...





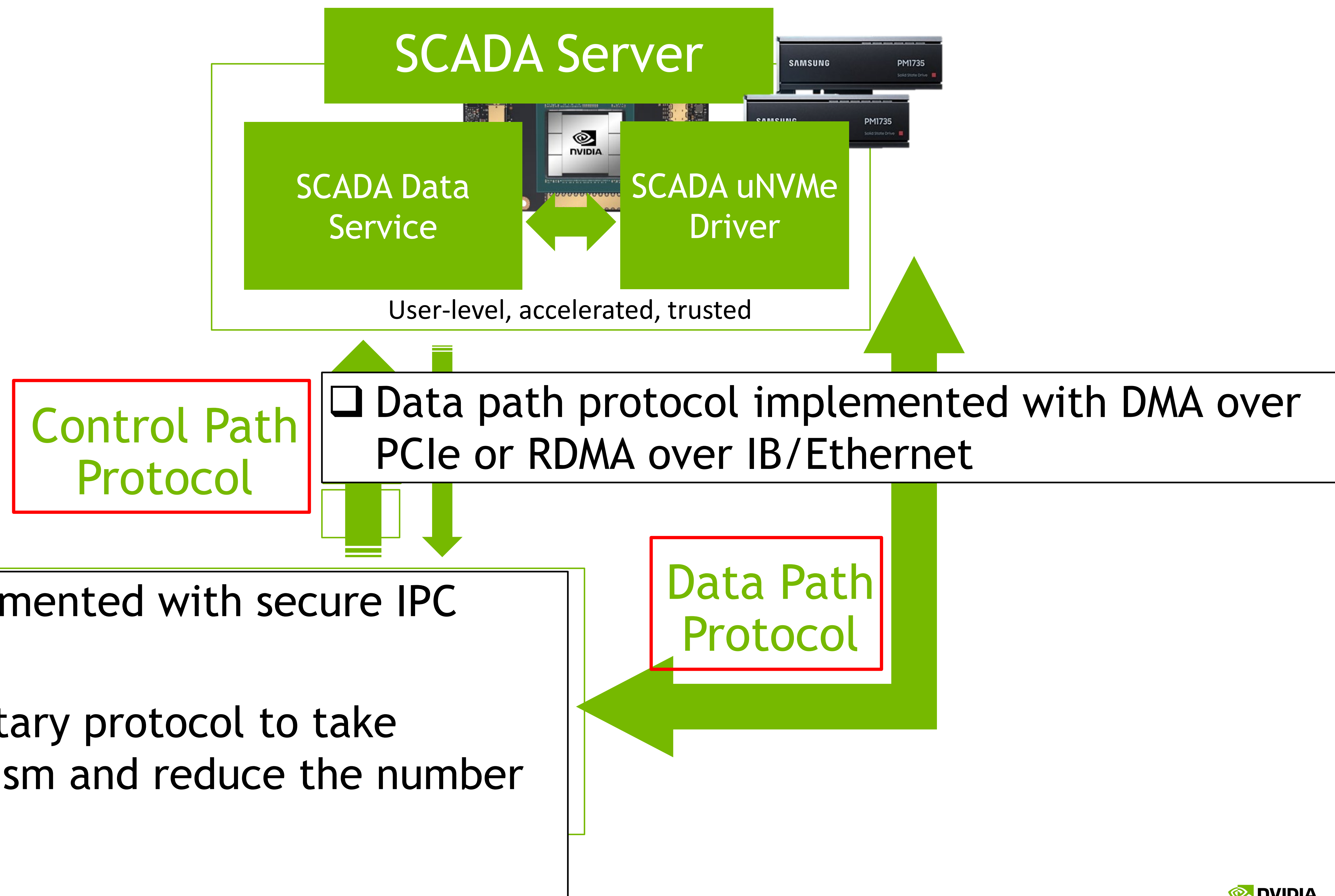
# SCADA – Scaled Accelerated Data Access

Software Architecture and Components



# SCADA – Scaled Accelerated Data Access

Software Architecture and Components





# H3 Platform Gen5 System

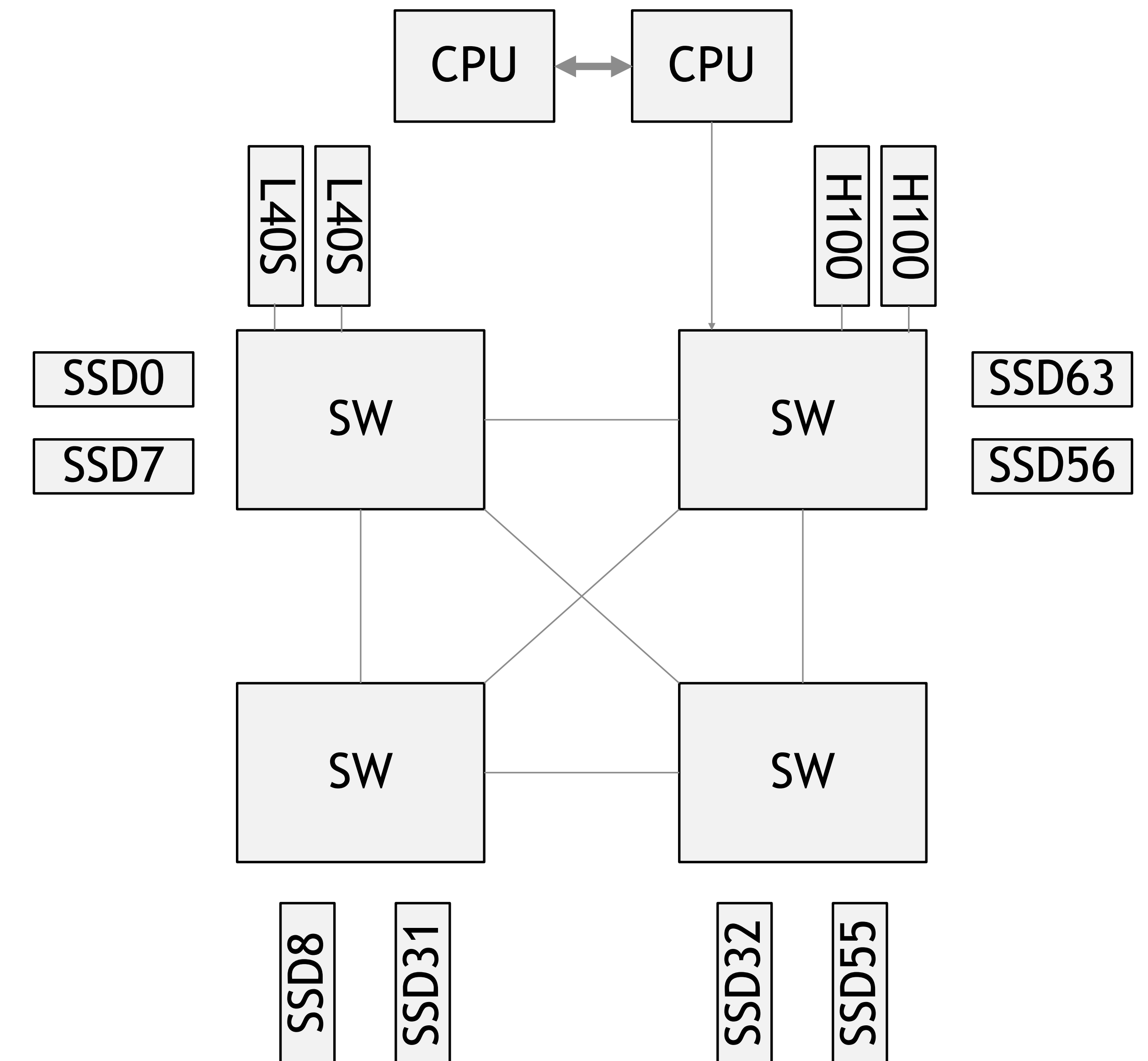


**CPU BOX**

**SSD BOX**

# H3 Platform Configuration

- INTEL(R) XEON(R) PLATINUM 8568Y
- 192 CPU cores
- 1 TB Memory
- 2 Micron SSD for OS.
- 2 Micron SSD for Data drives
- 64 Samsung E1.S PM9D3A SSDs for SCADA



Config as of today

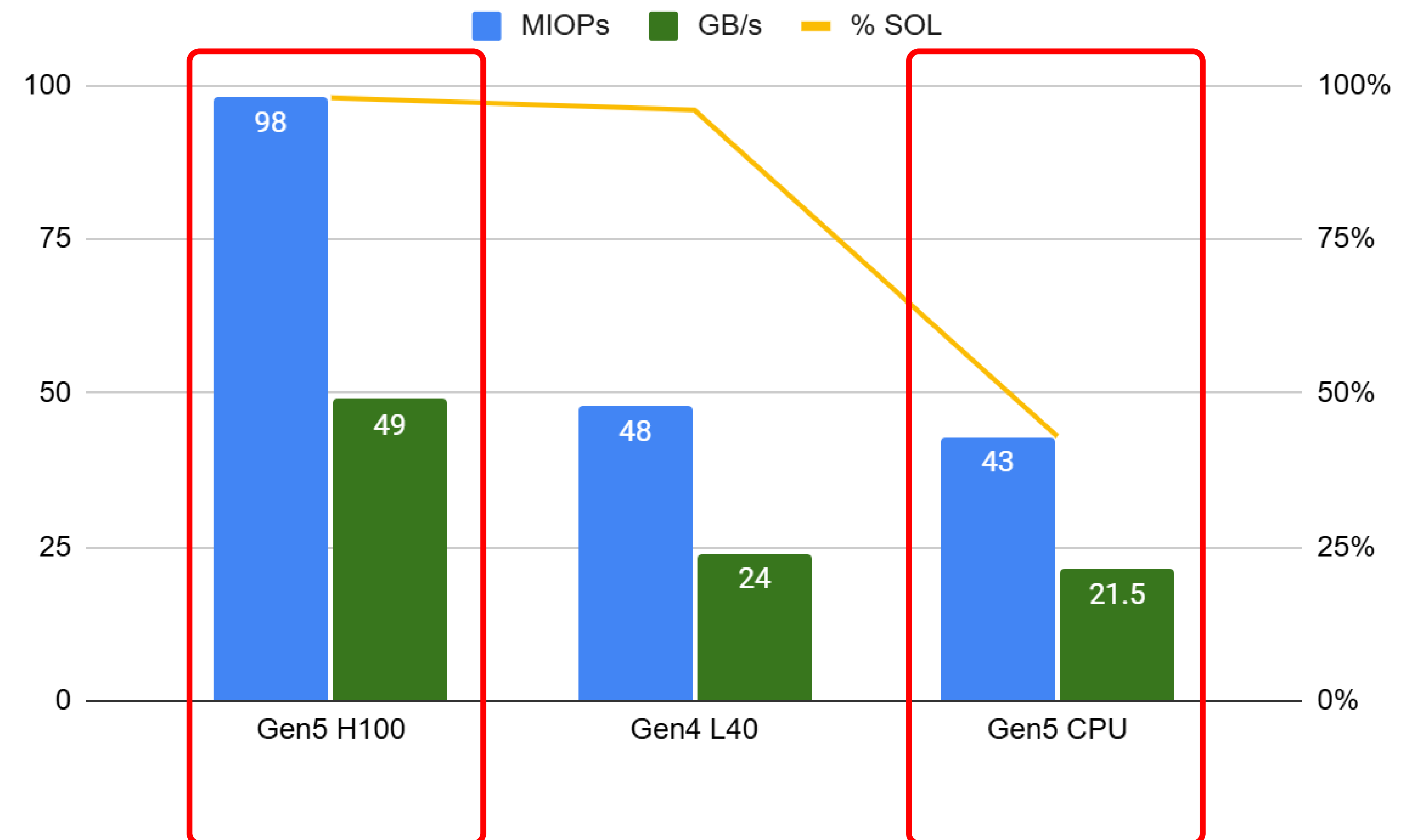


# Initial measurement summary

Implications for storage server: GPU vs. CPU (2x) to saturate IOPs, GB/s; PCIe switches in fabric node are efficient

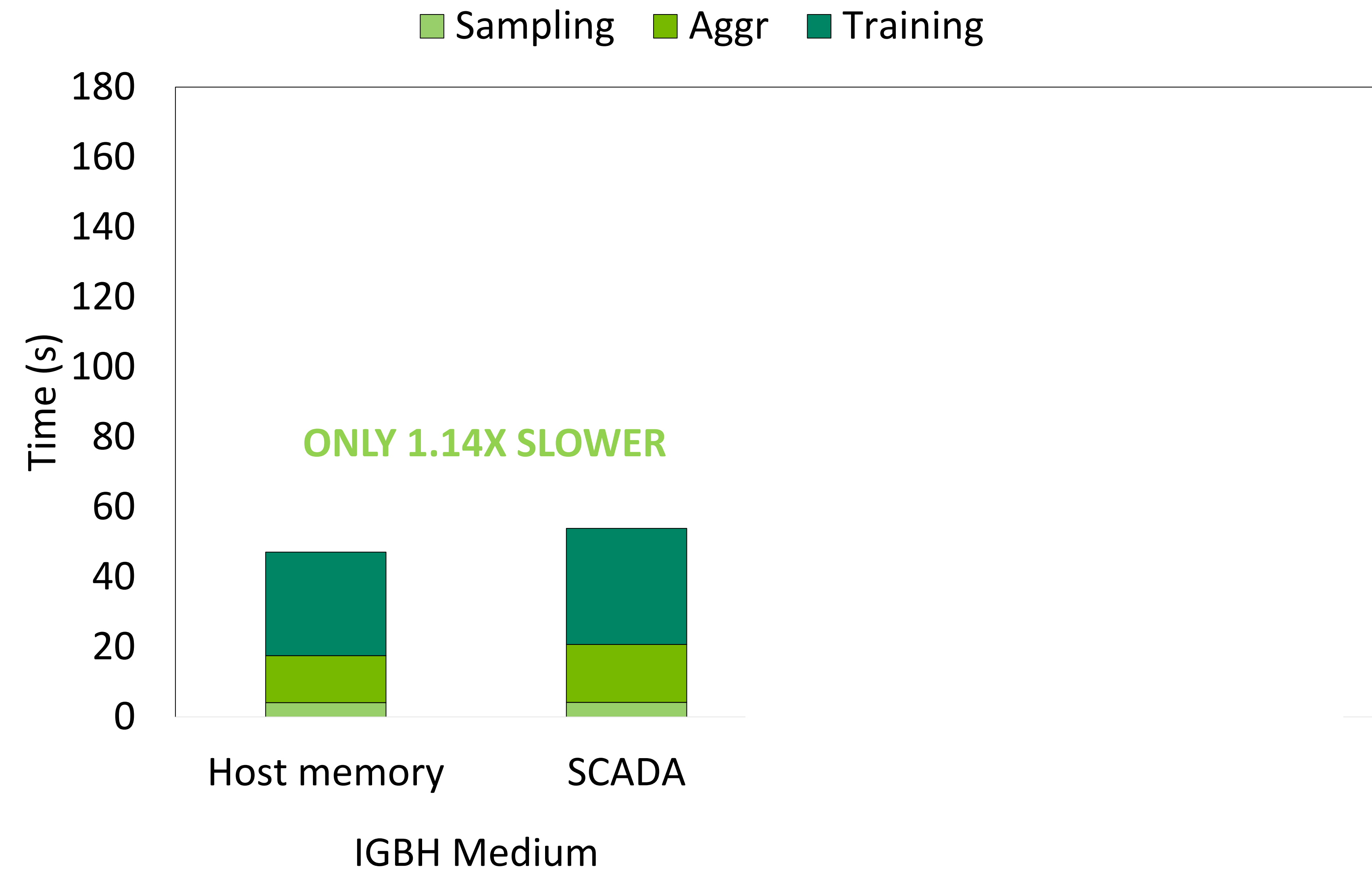
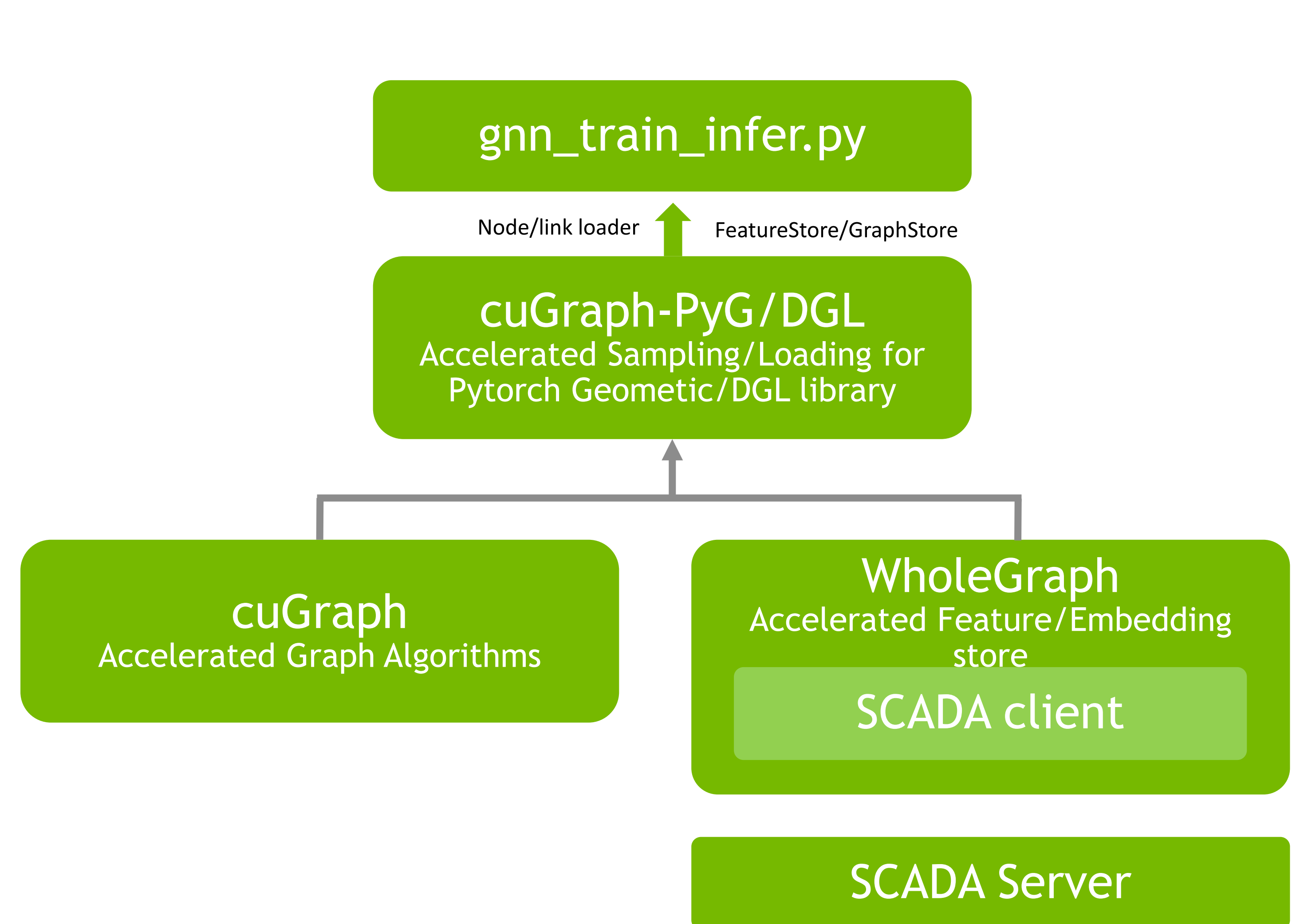
- Measurements from NVMe driver only
  - Full E2E results still being tuned to use more request buffers
  - Measurements made on H3 Storage-Next Experiment server
- H100 nearly saturates IOPs, GB/s
  - Broadcom PCIe switches in fabric mode very efficient
  - 96% of peak bandwidth and 512B IOPs
  - **Achieves that with 24% utilization, not yet tuned**
- Proof! CPU is insufficient to keep up
  - Gen5 GPU is >2x as effective as a Gen5 CPU
  - 1x Gen5 Intel Xeon Platinum, 48C 4GHz
- Modest GPU's perf isn't bad for its Gen
  - L40S is Gen4, there's no equivalent in Gen5
    - Based on Ada, not Hopper. Fewer SMs, lower HBM bandwidth
  - Cheaper GPU is 98% as efficient ( $48/98 \times 2$ ) for its generation

Compute Agent Comparison in SCADA Experiment Box



# SCADA GNN Training

Preliminary results using MLPerf GNN benchmark on IGBH dataset using Gen5 System

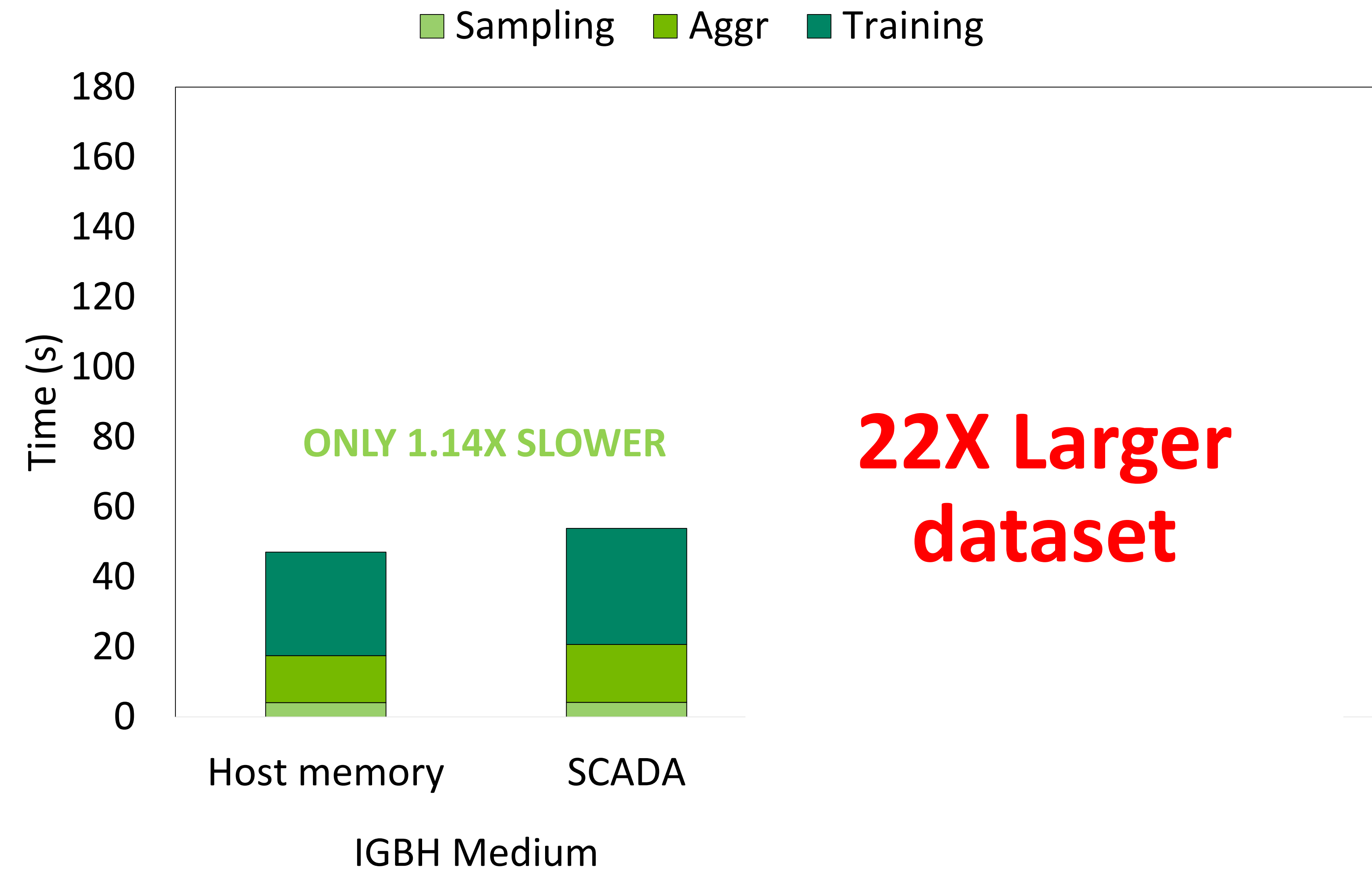
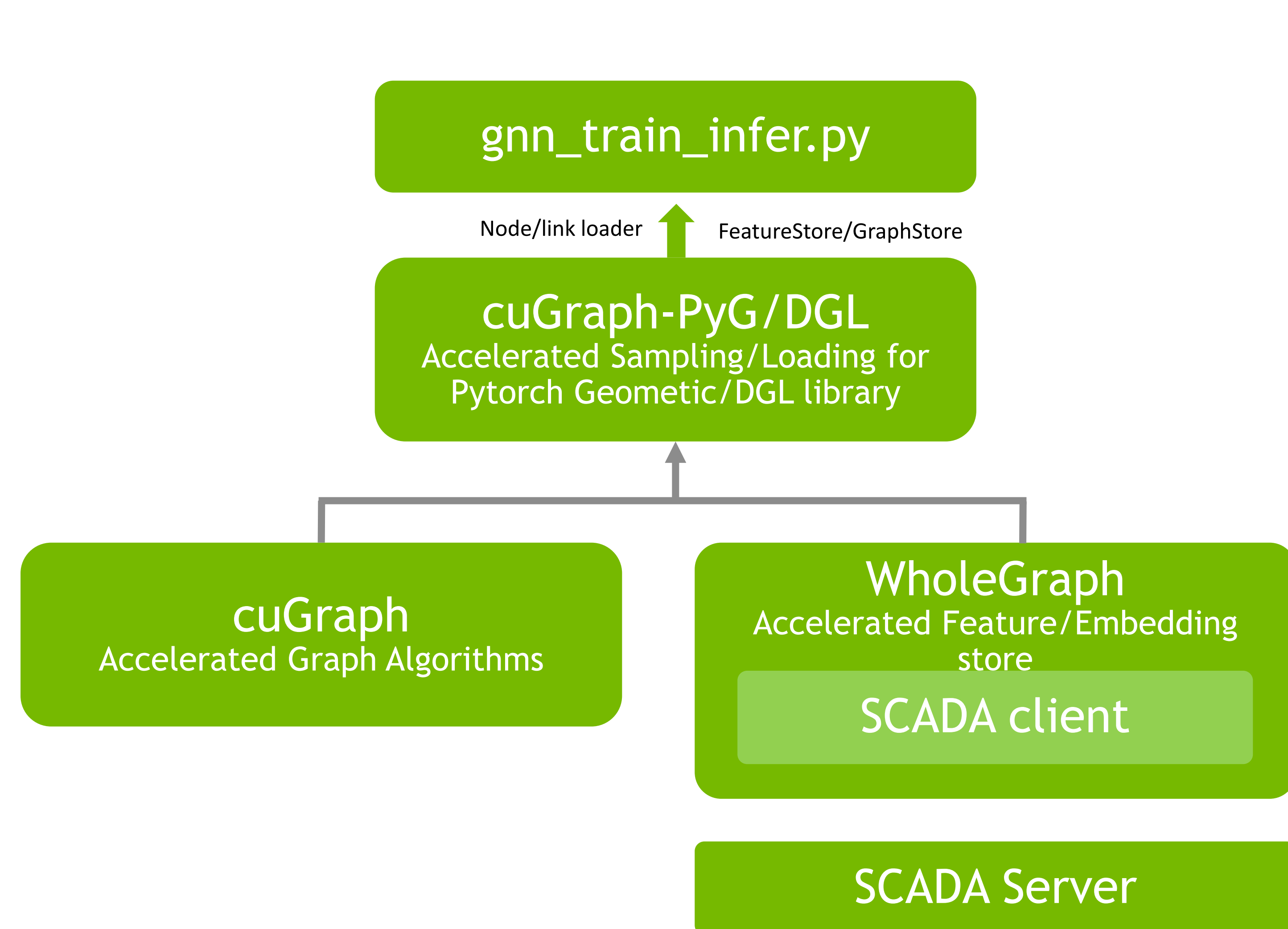


**SCADA DOES NOT USE CPU MEMORY IN THIS EXPERIMENT**

1 H100 GPU with 4 Samsung PM9D3A Gen5 drives  
IGBH-medium run with bz=2K, IGBH-Full with bz=1K, 4KB access.  
GNN RGAT model

# SCADA GNN Training

Preliminary results using MLPerf GNN benchmark on IGBH dataset using Gen5 System



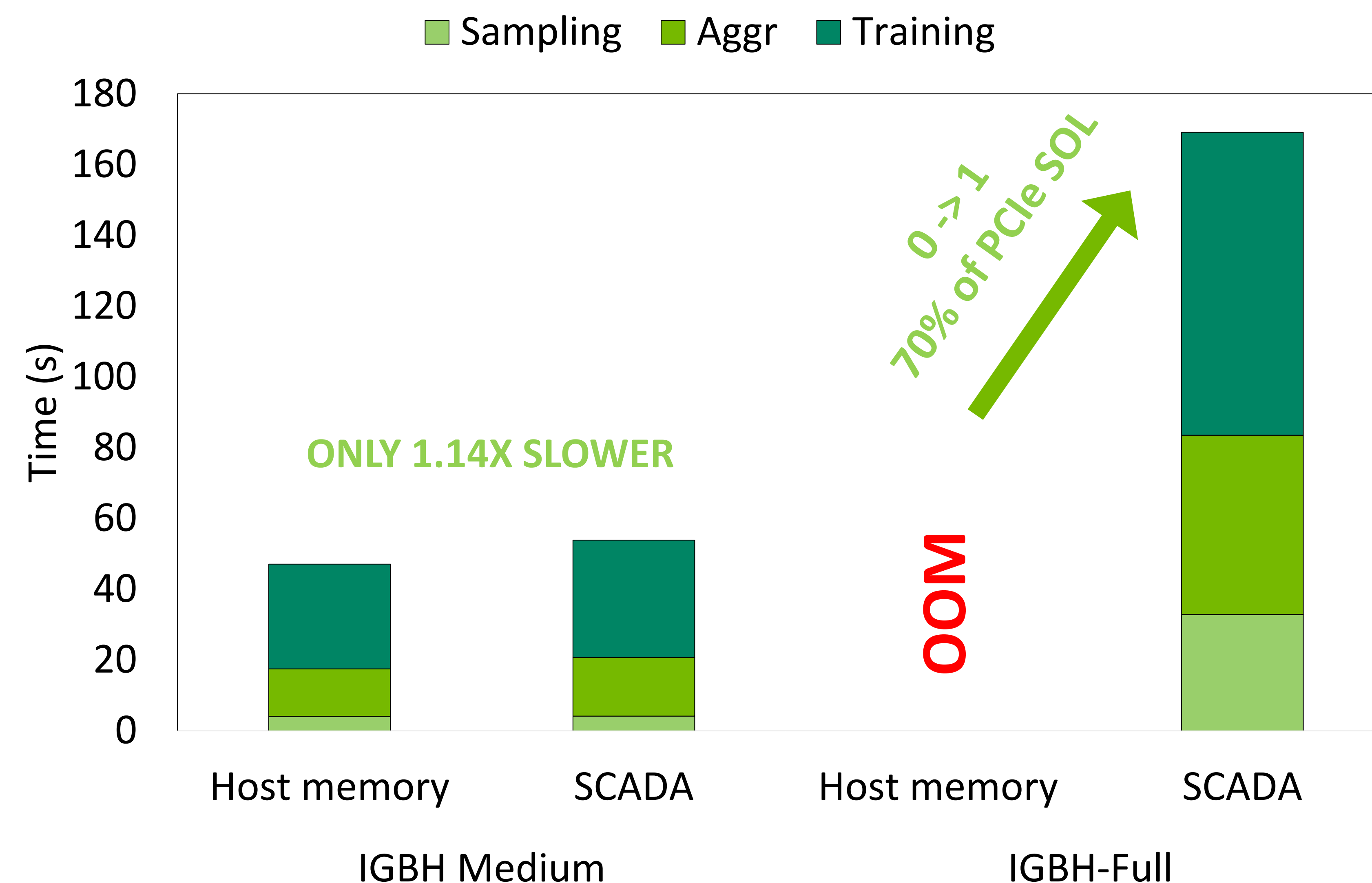
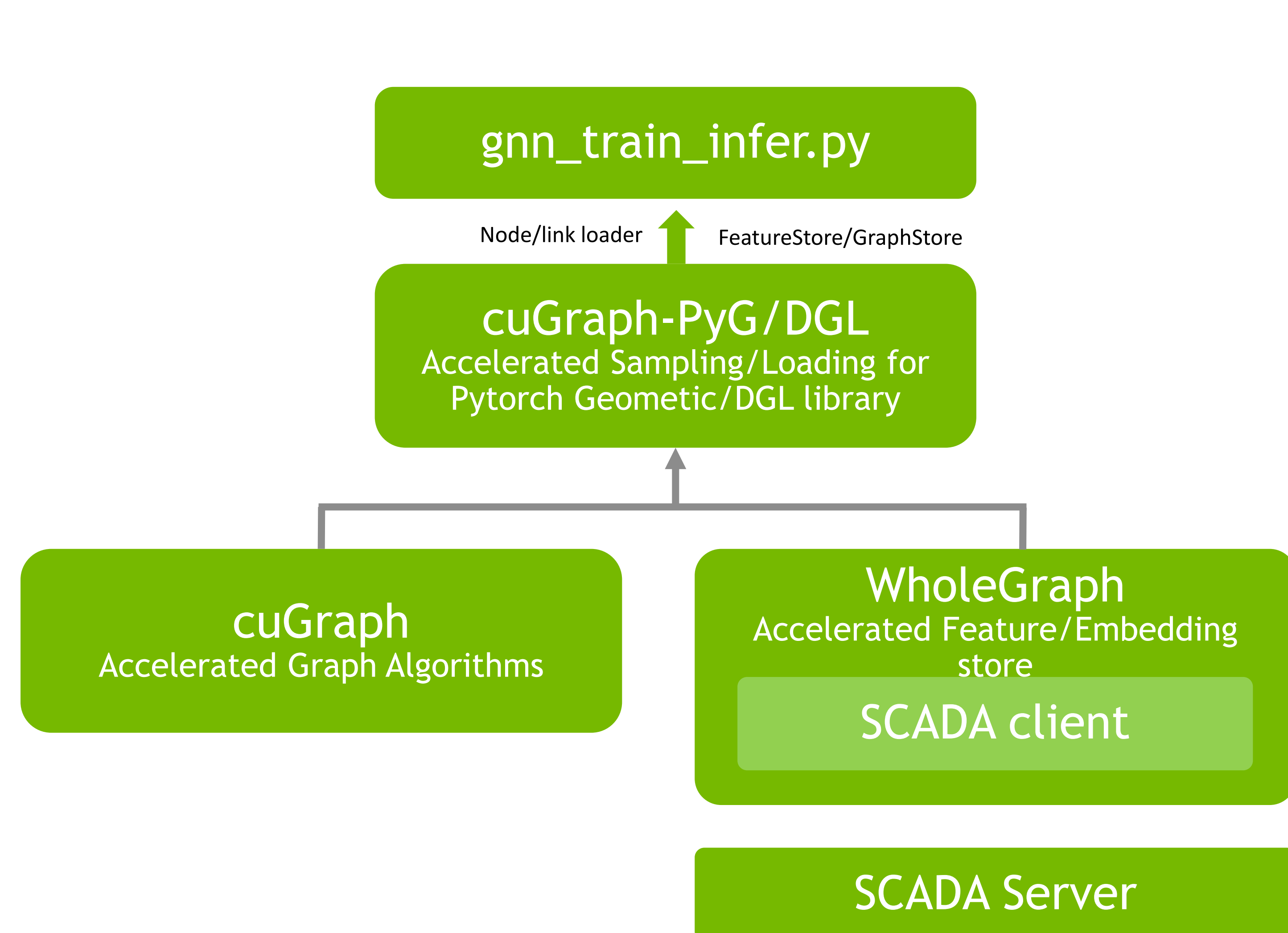
**SCADA DOES NOT USE CPU MEMORY IN THIS EXPERIMENT**

1 H100 GPU with 4 Samsung PM9D3A Gen5 drives  
IGBH-medium run with bz=2K, IGBH-Full with bz=1K, 4KB access  
GNN RGAT model



# SCADA GNN Training

Preliminary results using MLPerf GNN benchmark on IGBH dataset using Gen5 System

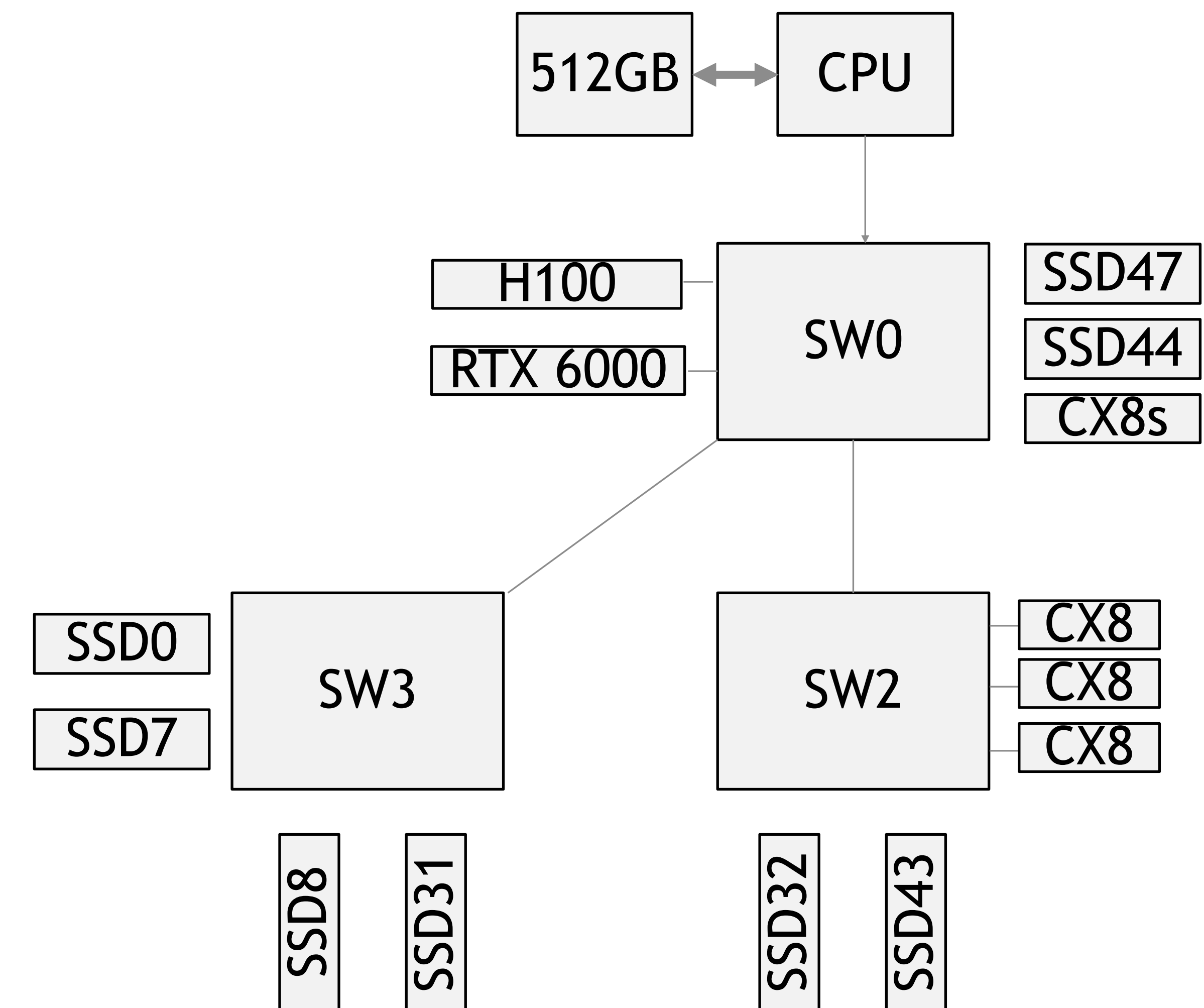
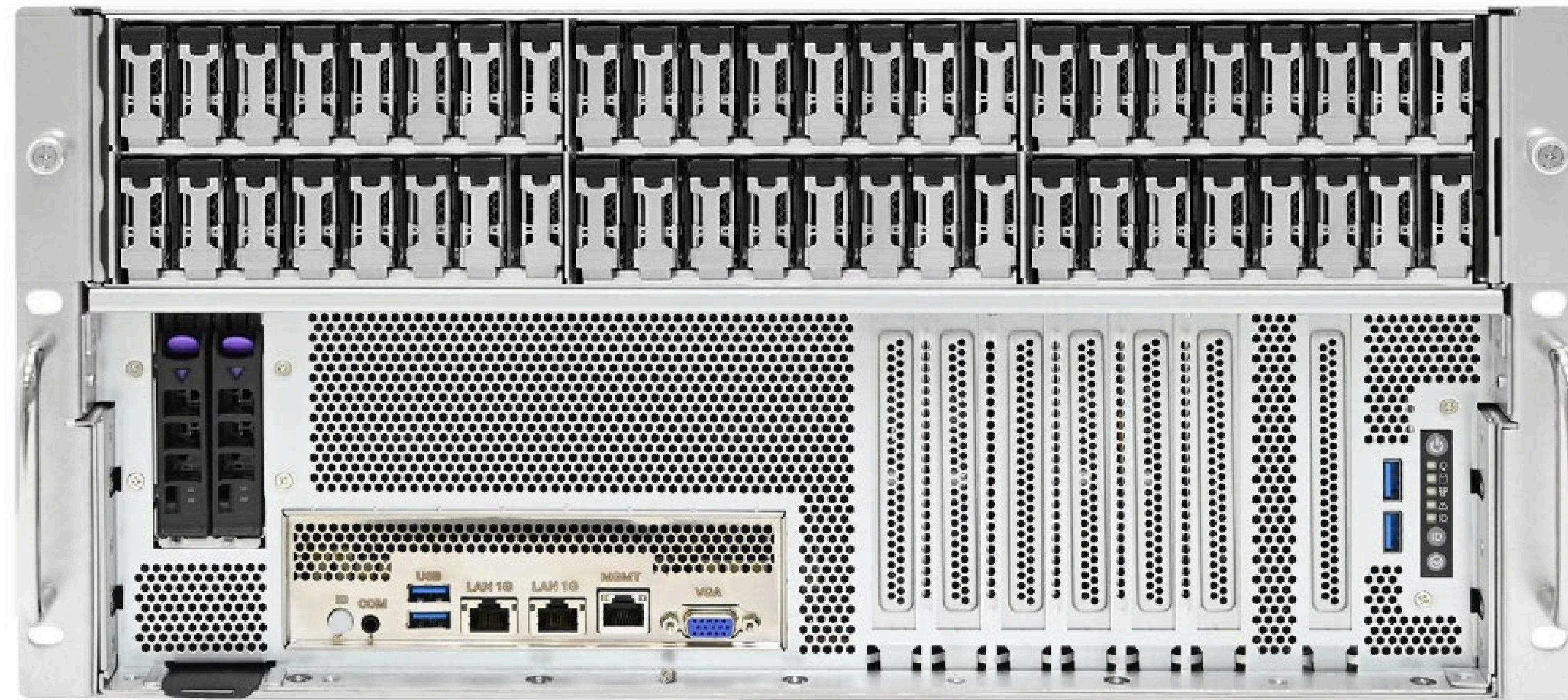


**SCADA DOES NOT USE CPU MEMORY IN THIS EXPERIMENT**

Single H100 GPU with 4 Samsung PM9D3A Gen5 drives  
 IGBH-medium run with bz=2K, IGBH-Full with bz=1K, 4KB access.  
 GNN RGAT model

# H3 Platform Gen6 System

48 E1.S SSD Platform with up to 4 GPUs and 6 CX8s

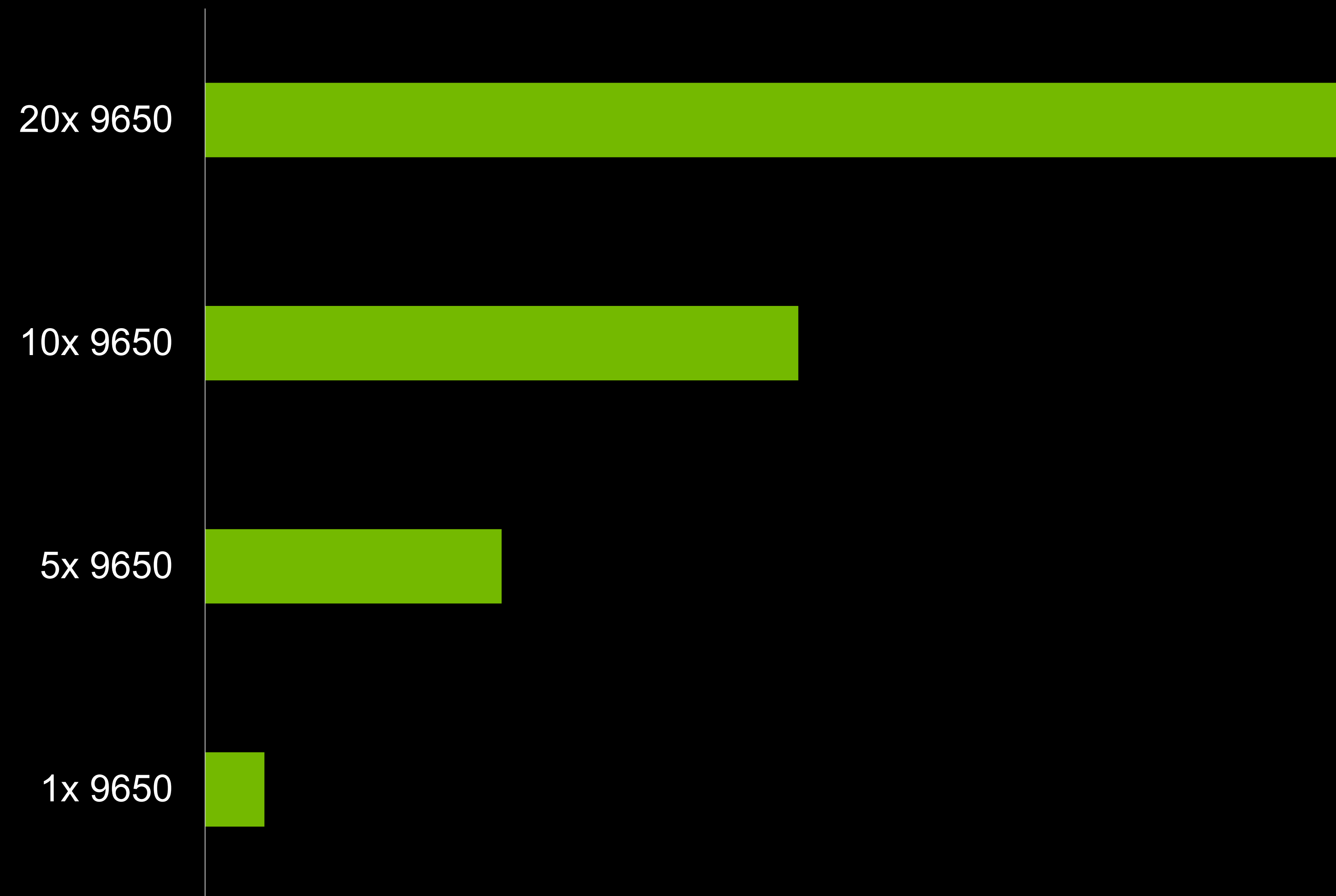


- Expected performance
  - 48 Micron Gen6 SSDs
  - $48 \text{ E1.S} * 5.4 \text{ M IOPS per SSD} = 260 \text{ M IOPS}$
  - Upto seq bandwidth – 600GBps

# Micron 9650 | NVIDIA SCADA

Early SCADA code shows strong performance and linear scaling on H3 Platform (Preliminary Results)

## ● NVIDIA SCADA (Preview Software)



**Early NVIDIA SCADA code drives impressive small block random read IOPs through 20 Micron 9650 Gen6 NVMe SSDs**

- **Linear performance scaling from 1 to 20 drives**

- **H3 Platform System:**

- Intel 8568Y+, 512GB DDR5
- 3x Broadcom 144 lane PCIe Gen6 A0 Switches
- 20x Micron 9650 Gen6 NVMe SSD, E1.S 7.68TB
- H100 NVL 96GB HBM3

- \* **Preliminary Results:**

- \* **Hardware & software stack tuning for ongoing**

SCADA test results collected on pre-production code.

System under test: H3 Platform PoC system, 1x Intel 8568Y+, 512GB DDR5, 3x Broadcom A0 PCIe Gen6 switches, 20x Micron 9650 E1.S 7.68TB, NVIDIA H100NVL-96GB PCIe Gen5x16, Workload is 512B random read initiated from H100 GPU.

Performance testing completed by Micron's Data Center Workload Engineering team.



# Collaborators

This is Marathon and not a sprint!

RAG, Fraud detection,  
ads, recsys, other apps...

Framework providers

CSPs, Hyperscaler and  
enterprise players

OS, hypervisors, driver and  
library providers

Storage Partners

OEM/ODM partners

SSD and DRAM media

Hyperscalers, OEMs

Storage  
Providers

NVIDIA

Media controllers  
Media

Artist's conception of a timeline

Gen5 Gen6 CSP feedback Cooling feedback

Biz/usage models  
Requirements Provide  
workloads Evaluation  
Characterization Refined  
requirements Interop  
requirements

Concept Kickoff Share  
ideas Refined  
usage Experiments  
infra Evaluation  
characterization Refined  
requirements

Kickoff  
fDec'24

GTC  
Mar'25

FMS  
Aug'25

SDC  
Sep'25

OCP  
Oct'25

Std?

F A D U

Graid Technology Inc.

KIOXIA

MARVELL

MICROCHIP

micron

PHISON

PLiOPS  
EXTREME DATA PROCESSOR

SAMSUNG

SANDISK

ScaleFlux

SiliconMotion

SK hynix

SmartIOPS

SOLIDIGM

Western Digital

AIC

ddn

DELL Technologies

H3

Hewlett Packard  
Enterprise

Hitachi Vantara

IBM

NetApp

PURESTORAGE

VAST

WEKA

## We have started! Engage with us!

Logos shown are owned by their respective companies.

NVIDIA



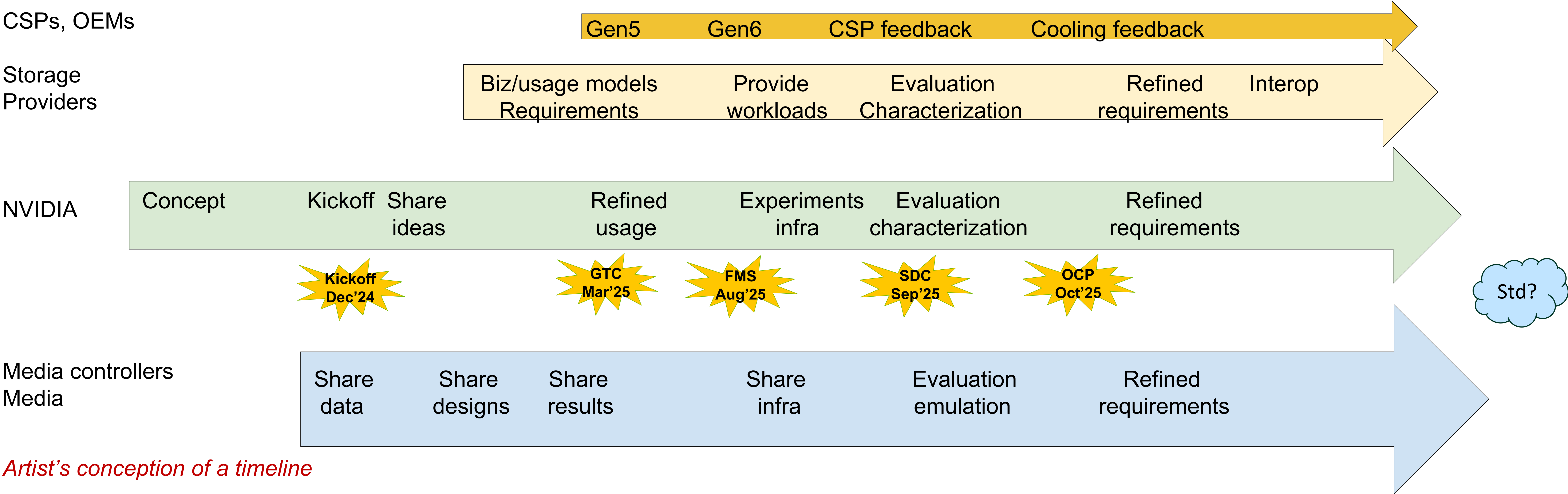






# The path to Storage-Next

This is Marathon and not a sprint



Artist's conception of a timeline

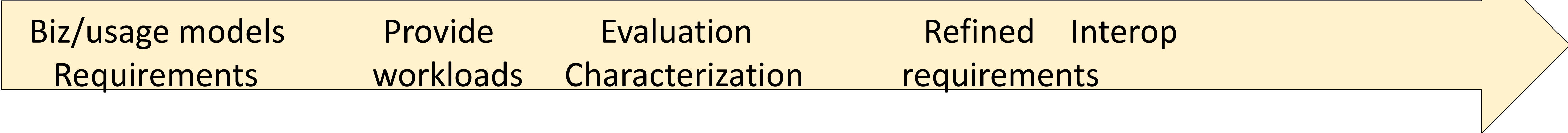




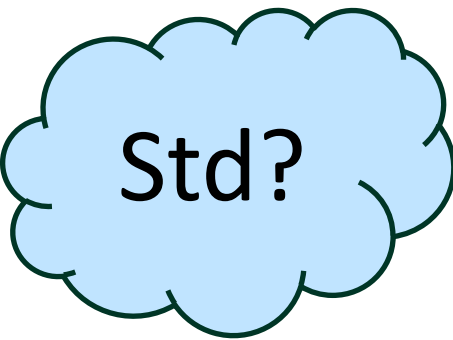
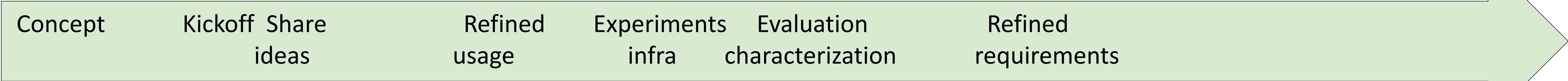
Hyperscalers, OEMs



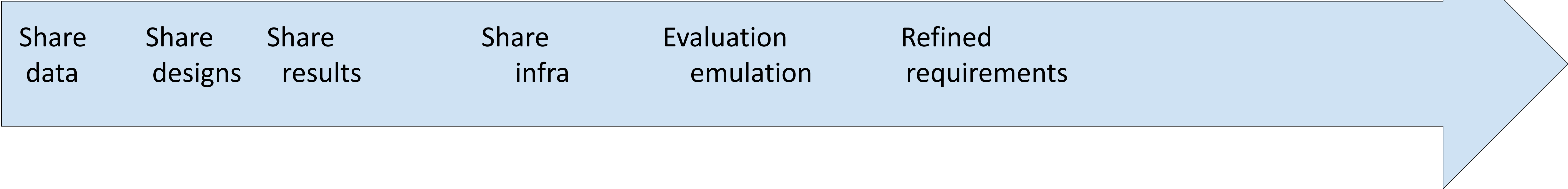
Storage Providers



NVIDIA



Media controllers  
Media



Artist's conception of a timeline



# Real-Time AI - Converged Data and Inference Applications

ms-level latency for fine-grained access to PB-scale data during inference

