

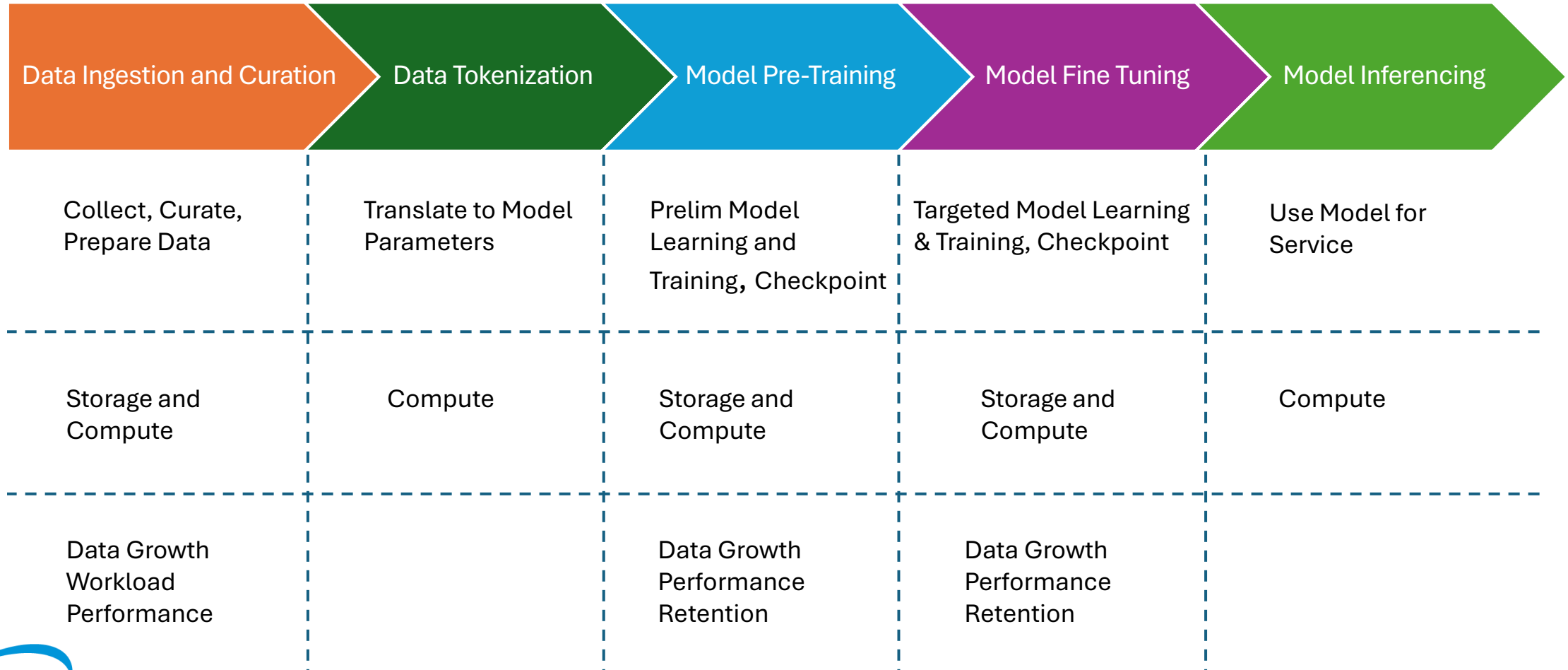
# Storage for AI

Presented By

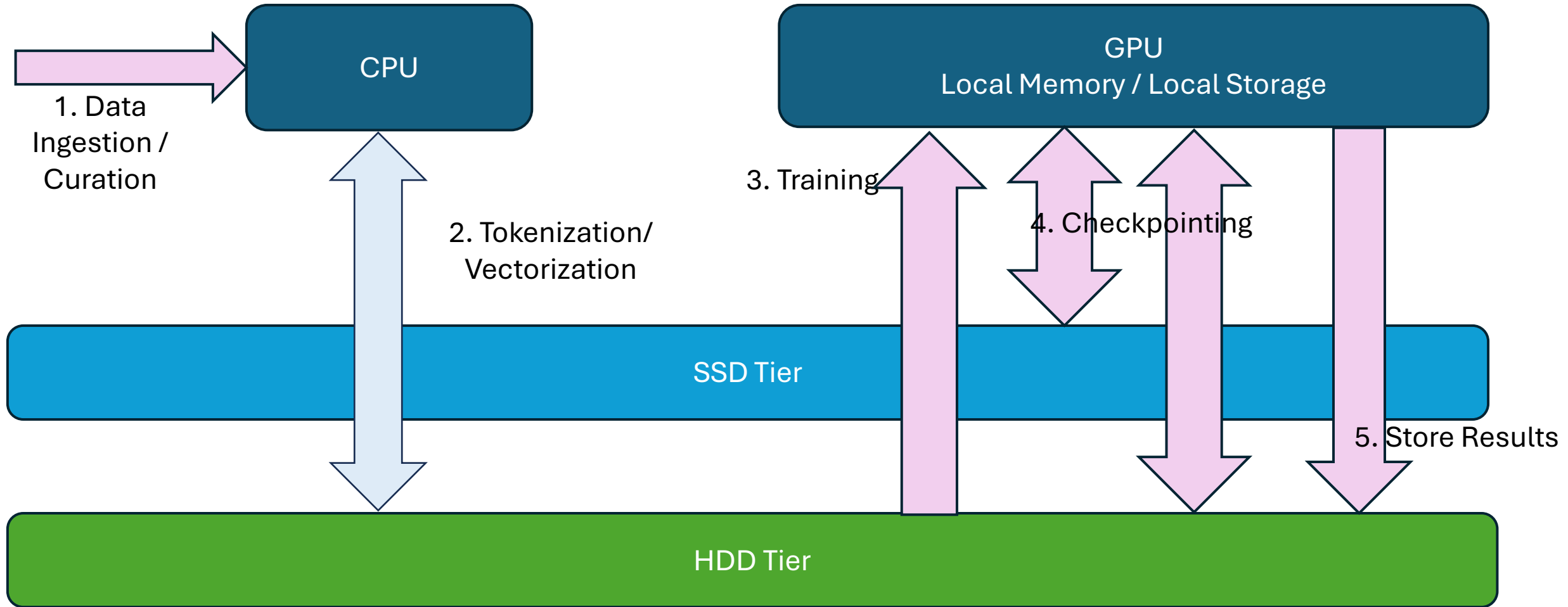
Shashidhar Joshi, Principal, Microsoft  
Swapna Yasarapu, Principal, Microsoft



# AI Model Training



# AI Data's Journey through Storage Tiers



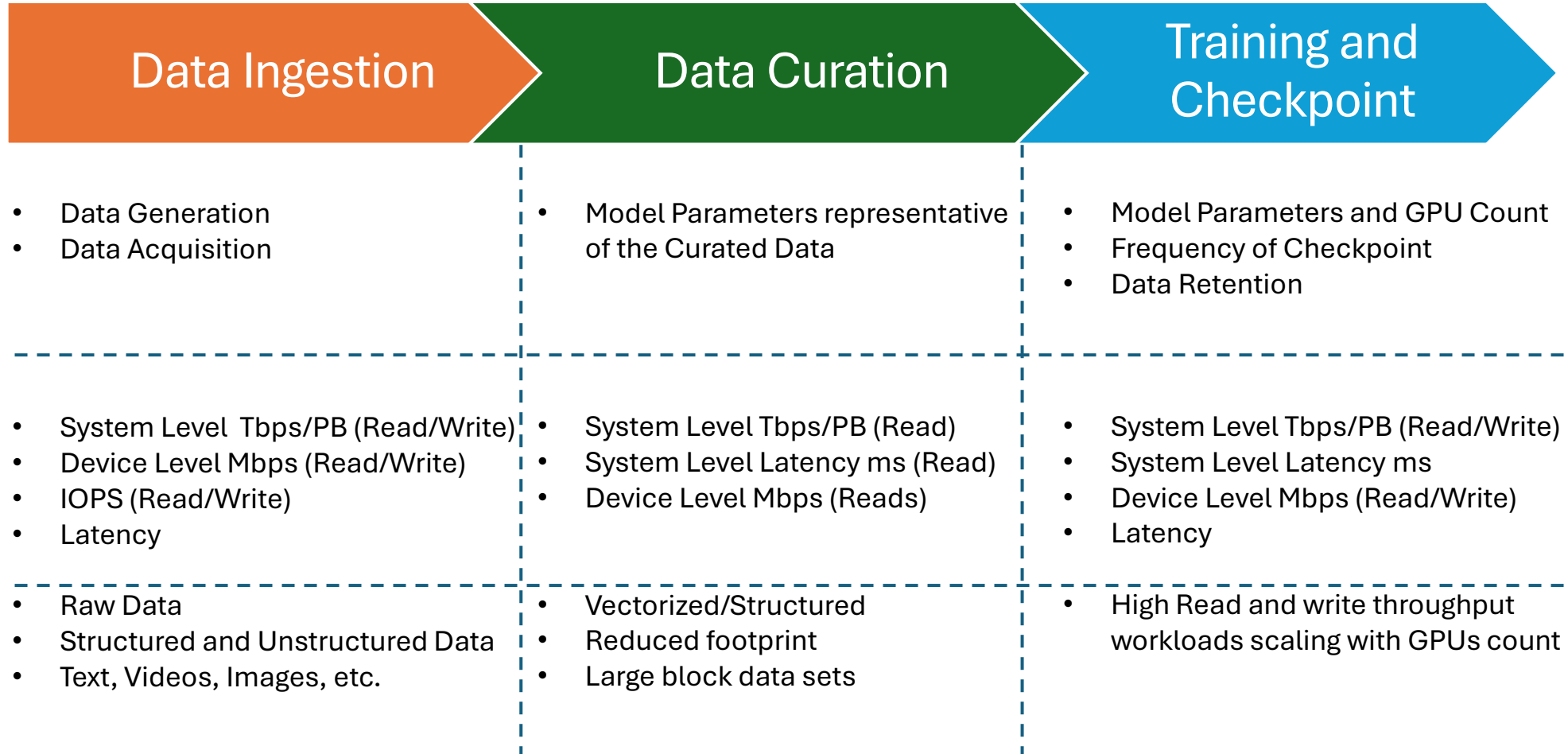
• Ingested and Curated Data is stored in our Blob Storage Service, mostly supported by HDDs

*the Future of Memory and Storage*

- Most Recent Checkpoint lands on SSD
- Previous Checkpoints land on HDDs for Longer Retention
- Scales with GPUs and Stop/Start Cycles

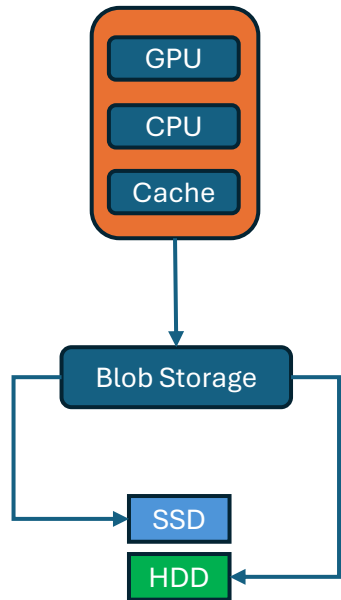


# AI Storage Demand and Performance

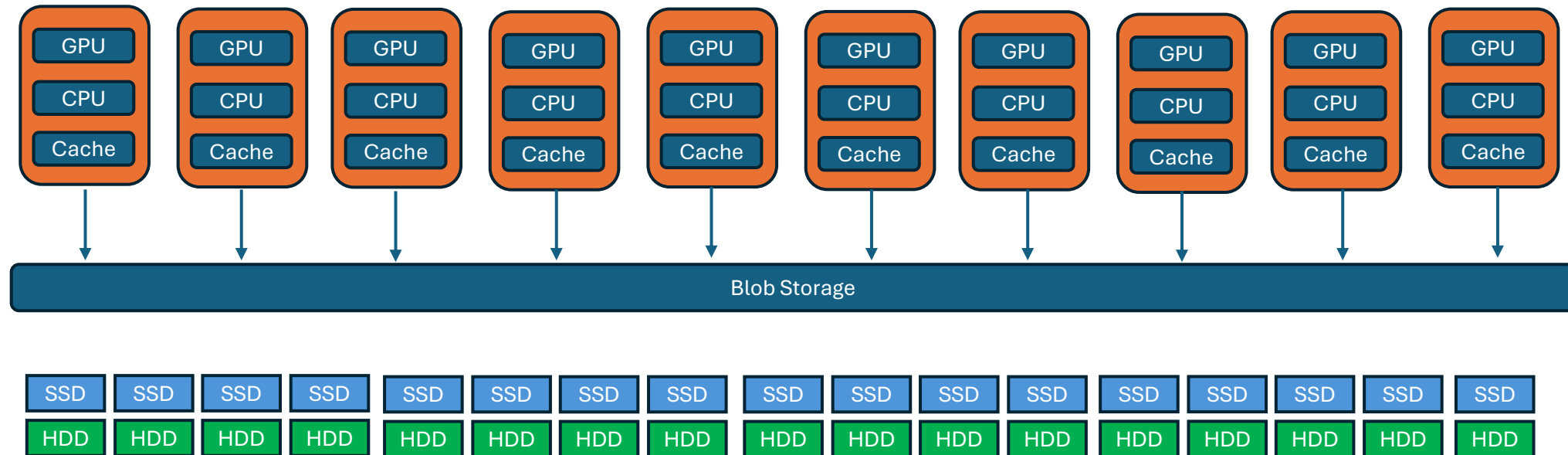


# Checkpoints: Driving Capacity and Thruput on Storage

A simple monolithic AI Training HW System



Scaling of GPUs as a function of Training Demand and Fetching of Checkpoints



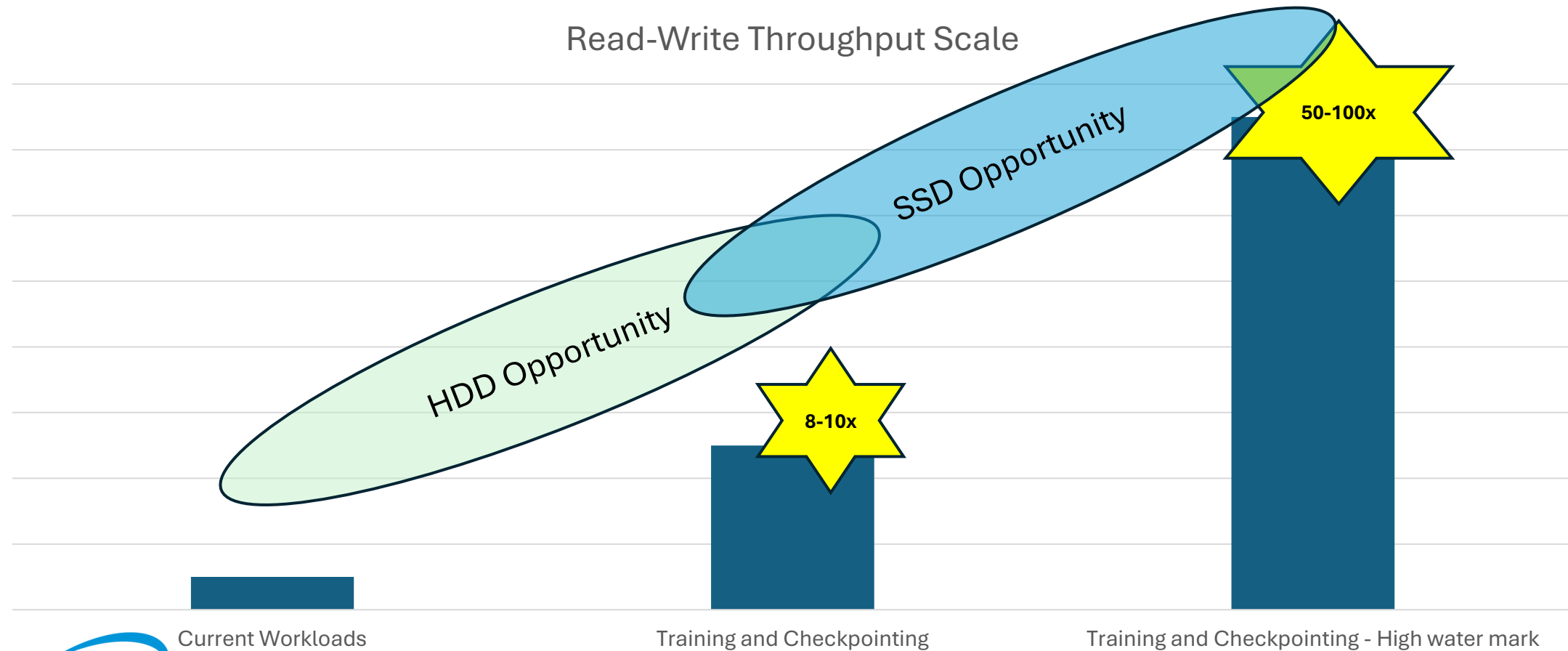
- Checkpoint Frequency: Model and Customer Unique. More Frequent Checkpoints lead to more Storage
- Most Recent Checkpoint Data: Most Recent copy available in SSD tier for low latency access
- Older Checkpoint Data: On HDD Blob Storage Tier, available but slow to access when needed
- GPU Scaling: All at once access from Blob Storage Tier drives high thruput



the Future of Memory and Storage



# AI Workload Implications: MB/s/TB



# Media Technologies - Call to Action

## *HDDs*

- Power Efficiency Management
- Scale Throughput Density with HDD Capacity

## *SSDs/Flash*

- Density Scaling
- Drive TCO Improvement while delivering on the throughput-Density-Power
- Step Function Improvement in Bit Cost

# Questions?