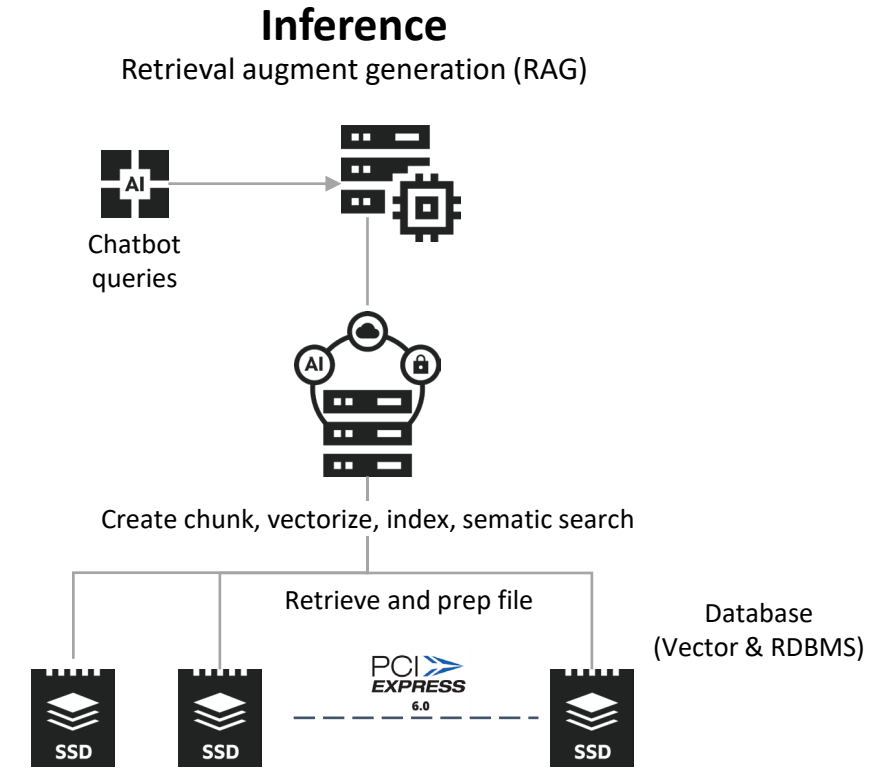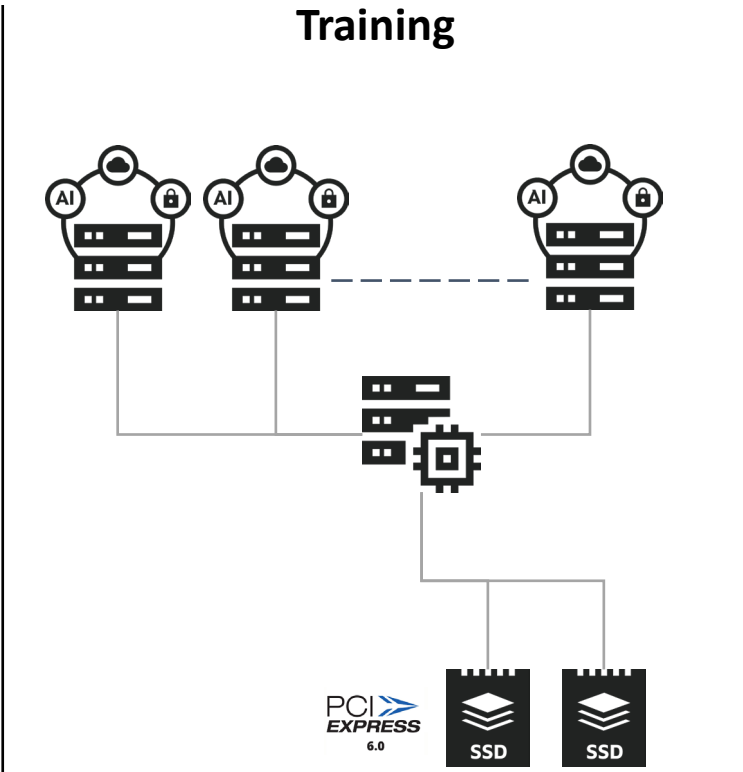# Flash Controller for the AI era
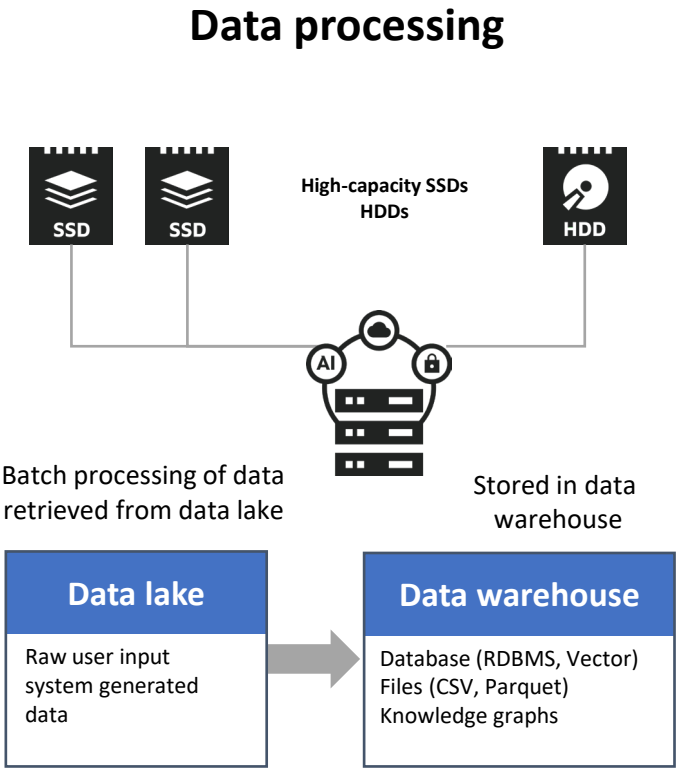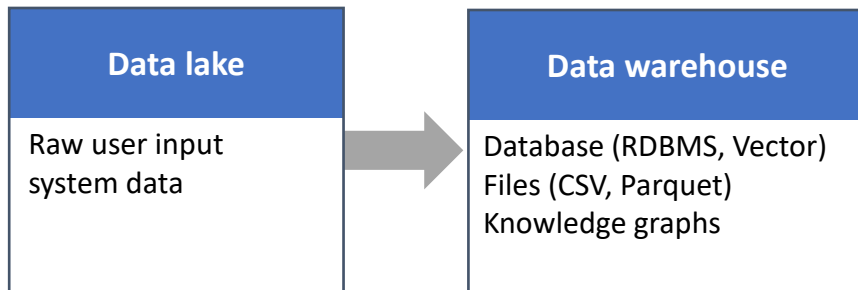
Vasanthi Jagatha, Senior Manager, Marvell

# AI workloads using storage

## Data processing

High-capacity SSDs
HDDs

Batch processing of data
retrieved from data lake

Stored in data
warehouse

| Data lake | Data warehouse |
|---|---|
| Raw user input system generated data | Database (RDBMS, Vector) Files (CSV, Parquet) Knowledge graphs |

## Training

PCI EXPRESS 6.0

## Inference
Retrieval augment generation (RAG)

Chatbot
queries

Create chunk, vectorize, index, sematic search

Retrieve and prep file

Database
(Vector & RDBMS)

PCI EXPRESS 6.0

**Storage critical across data processing, training and inference workloads**

# Data processing



High-capacity SSDs HDDs

Batch processing of data retrieved from data lake

Stored in data warehouse

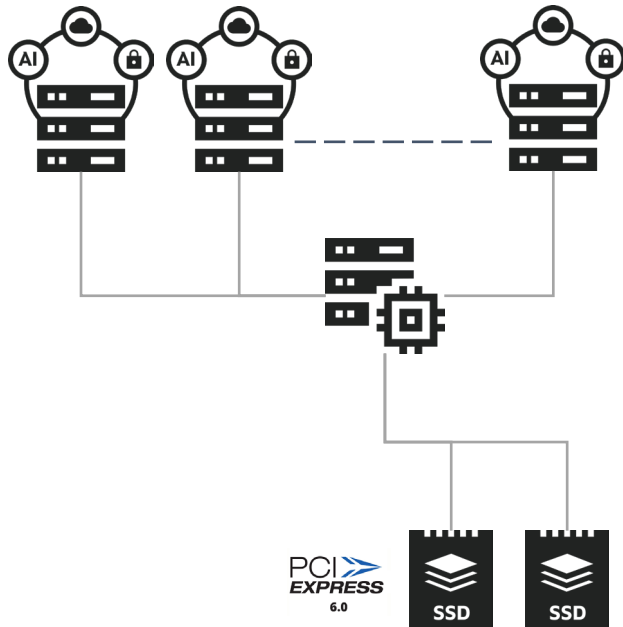| Data lake | Data warehouse |
|---|---|
| Raw user input system data | Database (RDBMS, Vector) Files (CSV, Parquet) Knowledge graphs |

- Dominates AI/ML development lifecycle
- Input data fidelity has outsize impact on resulting model performance
- Large data-sets typically stored in
  - Data lakes (unstructured)
  - Database, CSV, Parquet, JSON (structured)
- Bursty reads, write heavy workload

**Key Storage considerations- capacity, data reliability, and performance**
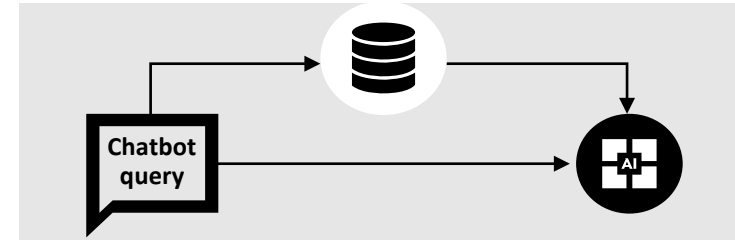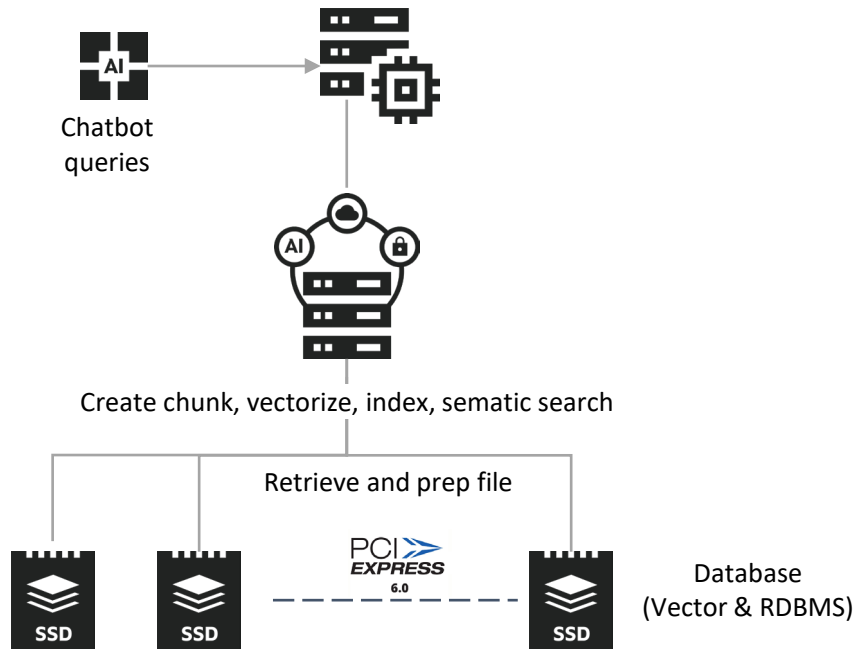
# Training



- Fast data retrieval -Reduce GPU idling
- Regular checkpointing  for observability
- Mixed reads and bursty writes

**High performance fast data retrieval and checkpointing**

# Inference – Retrieval Augmented Generation (RAG)



Chatbot queries

Create chunk, vectorize, index, sematic search

Retrieve and prep file

PCI EXPRESS 6.0
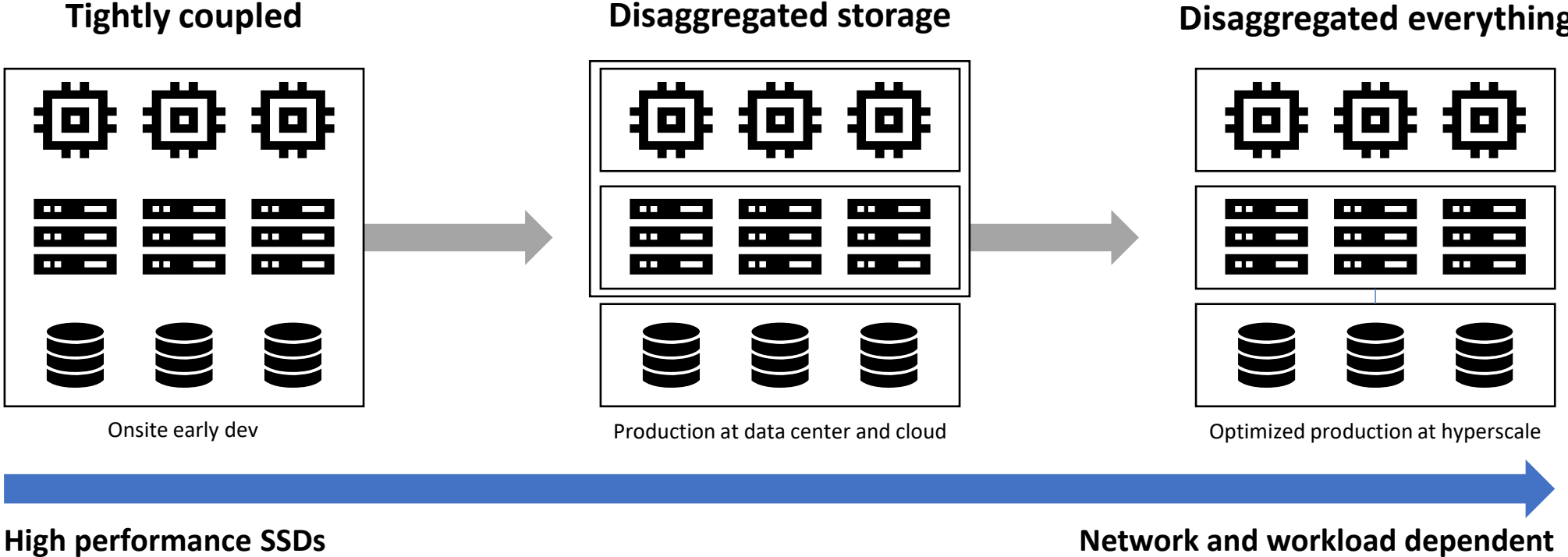
SSD   SSD   SSD

Database (Vector & RDBMS)

**RAG enhances LLMs by integrating external data retrieval**

- I/O Intensive (data prep, high user traffic)
- Large capacity  vector storage

**High performance, large capacity  for large scale RAG**

# AI infrastructure

**Tightly coupled**

**Disaggregated storage**

**Disaggregated everything**

Onsite early dev

Production at data center and cloud

Optimized production at hyperscale

**High performance SSDs**

**Network and workload dependent**

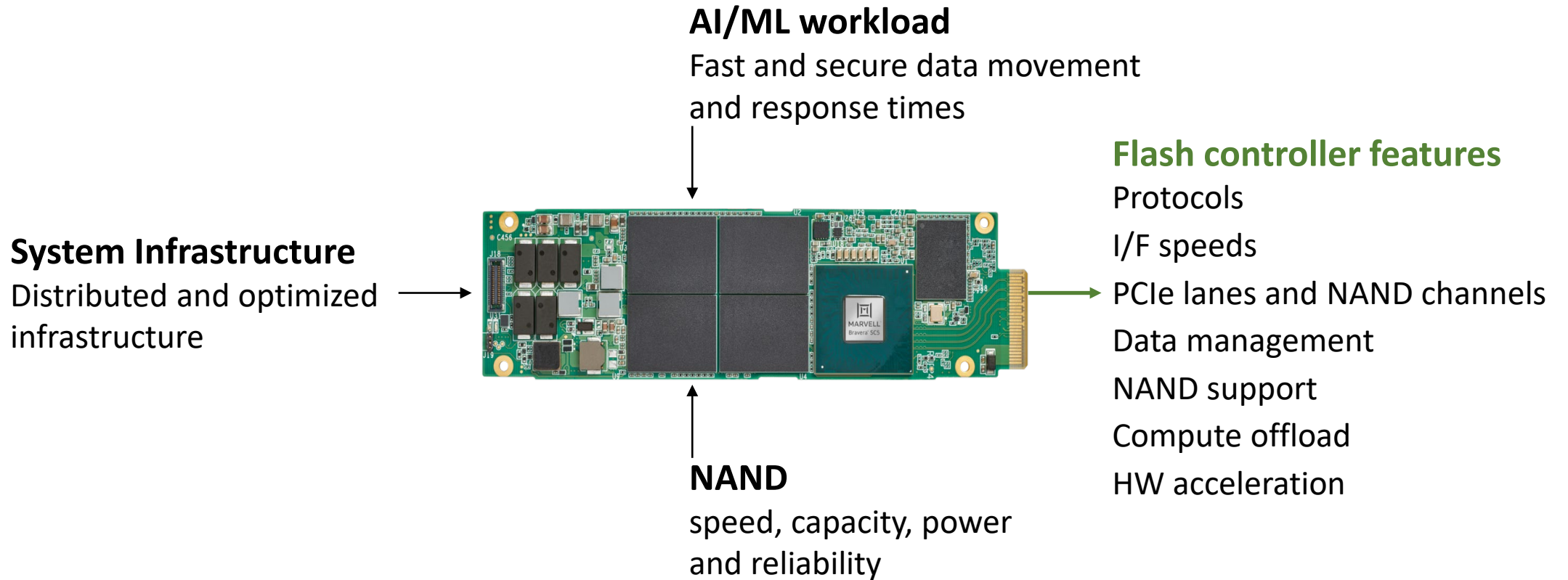**AI infrastructure needs evolves with AI/ML development phase - impacts storage requirements**

# Flash controller design



**AI/ML workload**
Fast and secure data movement
and response times

**System Infrastructure**
Distributed and optimized
infrastructure

**NAND**
speed, capacity, power
and reliability

**Flash controller features**
Protocols
I/F speeds
PCIe lanes and NAND channels
Data management
NAND support
Compute offload
HW acceleration

**Flash controller requirements driven by AI/ML workloads, system, NAND decisions**

# March to Gen6 SSDs



**28GBps**
**NAND 4800MT/s,1-2Tb**
**256TB SSD**
**LDPC+++**

**PCIe 6.0**

**14GBps**
**NAND 2400MT/s,512-1Tb**
**128TB SSD**
**LDPC++**

**PCIe 5.0**

**7.4GBps**
**NAND 1600MT/s,512Gb**
**64TB SSD**
**LDPC+**

**PCIe 4.0**

**3.2GBps**
**NAND 800MT/s,256Gb**
**16TB SSD**
**LDPC**

**PCIe 3.0**

Performance per watt

Capacity

Reliability

**Critical Flash Controller AI workload needs**
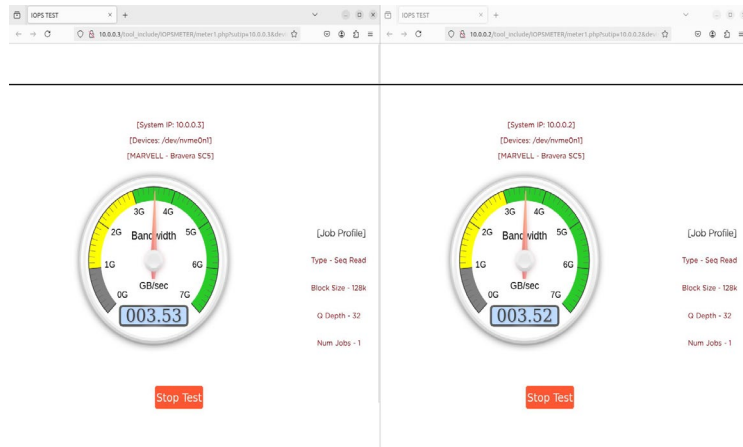
**Gen6 closes the gap for AI workload needs**

# Visit Marvell booth # 1046

**Dual Port Demo**          **Accelerated Storage for RAG and AI Inference**          **FDP QLC with PCIe 5.0 NVMe™ SSD**

**Customer Demo**



**Powered by**