# Today's Presenters



**Moderator:**
Chip Maurer
Senior Principal Engineer
Dell Technologies



**Presenter:**
Andy Longworth
Solution Architect
HPE



**Presenter:**
Vincent Hsu
VP, IBM Fellow, and CTO for
Storage and Software Defined
Infrastructure
IBM

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# SNIA-at-a-Glance

**180**
industry leading
organizations

**2,500**
active contributing
members

**50,000**
IT end users & storage
pros worldwide

Learn more: **snia.org/technical**   🐦 **@SNIA**

SNIA CSTI | CLOUD STORAGE TECHNOLOGIES

# SNIA Legal Notice

The material contained in this presentation is copyrighted by the SNIA unless otherwise noted.

Member companies and individual members may use this material in presentations and literature under the following conditions:

> Any slide or slides used must be reproduced in their entirety without modification
>
> The SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.

This presentation is a project of the SNIA.

Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.

The information presented herein represents the author's personal opinion and current understanding of the relevant issues involved. The author, the presenter, and the SNIA do not assume any responsibility or liability for damages arising out of any reliance on or use of this information.

NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.

SNIA.
CSTI | CLOUD STORAGE TECHNOLOGIES

# Agenda

- History of Big Data
- Current state
- Modernization challenges
- Evolving workloads, processing outside of data center
- Look towards the future

SNIA.
CSTI | CLOUD STORAGE
TECHNOLOGIES

# History of Big Data

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# History of Big Data

Enterprise Big Data Framework (https://www.bigdataframework.org/short-history-of-big-data/)

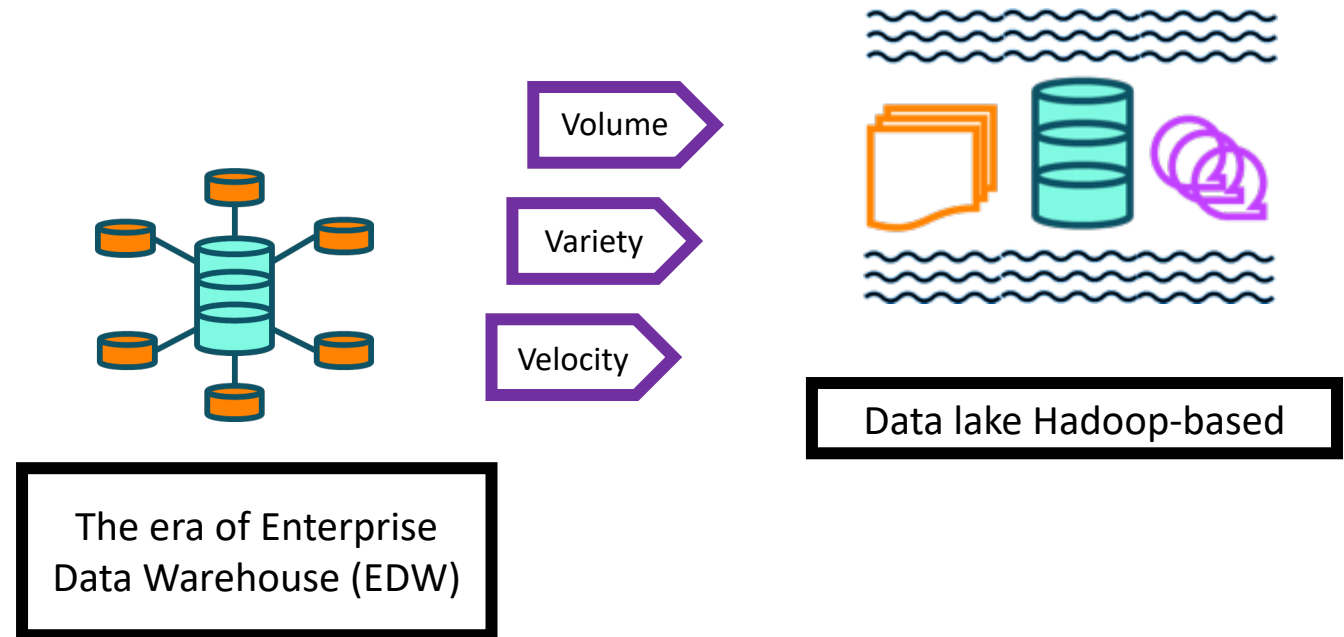| BIG DATA PHASE 1 | BIG DATA PHASE 2 | BIG DATA PHASE 3 |
|---|---|---|
| Period: 1970-2000 | Period: 2000-2010 | Period: 2010-present |
| DBMS-based, structured content: <br>• RDBMS & data warehousing <br>• Extract Transfer Load <br>• Online Analytical Processing <br>• Dashboards & scorecards <br>• Data mining & statistical analysis | Web-based, unstructured content <br>• Information retrieval and extraction <br>• Opinion mining <br>• Question answering <br>• Web analytics and web intelligence <br>• Social media analytics <br>• Social network analysis <br>• Spatial-temporal analysis | Mobile and sensor-based content <br>• Location-aware analysis <br>• Person-centered analysis <br>• Context-relevant analysis <br>• Mobile visualization <br>• Human-Computer-Interaction |

**Block, NFS, POSIX**          **HDFS**                              **Unified storage (Object and File)**

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Current State

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# Evolution of Data Platforms

- **Everything from Enterprise Data Warehouses to Hadoop based Data Lakes**

- **No one size fits all**

- **Emergence of cloud services**

- **Picking the right system for the right workload**
  - Structured vs Unstructured
  - Batch vs Real-time
  - On-premises vs Cloud

Volume

Variety

Velocity

The era of Enterprise Data Warehouse (EDW)

Data lake Hadoop-based

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

10

# The Five Vs of Big Data

- ## Started as 3 Vs
  - Volume: the huge amount of data that is produced every day
  - Variety: diversity of data, both types and sources
  - Velocity: the speed with which the data is generated

- ## Additional Vs
  - Veracity: is the authenticity and credibility of data
  - Value: transforming data into value for the business

SNIA CSTI | CLOUD STORAGE TECHNOLOGIES

# Is Hadoop Dead?

**For**

- Cost: running on commodity hardware
- Batch analytics
- Availability through fault tolerance
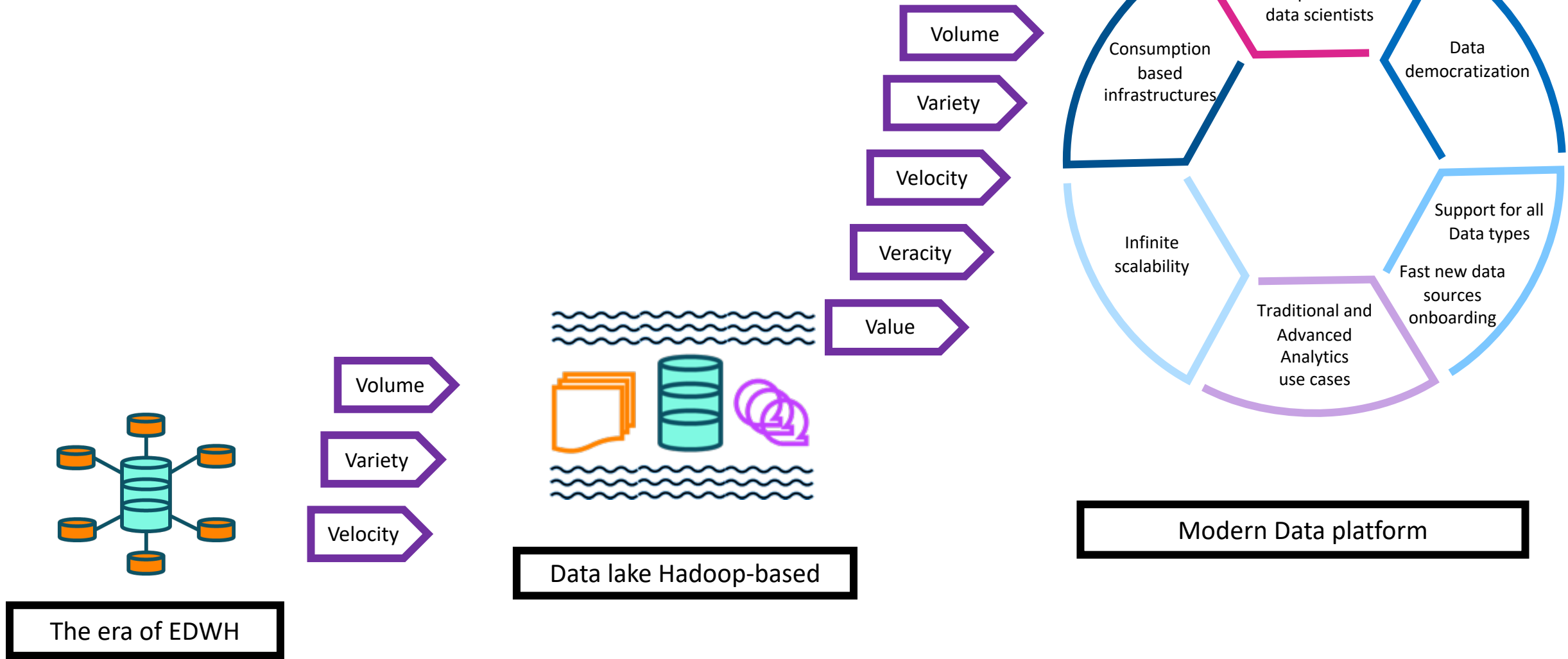- Spark on Hadoop

**Against**

- Inefficient for small datasets
- Real-time analytics
- Cloud alternatives
- Lack of integration with cloud services such as S3

- Merge of Hortonworks and Cloudera
- Cloudera Enterprise 6.2 & 6.3 EOL March 2022
- Hortonworks Data Platform 3.1 EOL December 2021

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Modernization Challenges

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# Evolution of Data Platforms

Volume

Variety

Velocity

Veracity

Value

Volume

Variety

Velocity

Volume

Variety

**Data lake Hadoop-based**

**The era of EDWH**

Enabled and empowered data scientists

Consumption based infrastructures

Data democratization

Infinite scalability

Support for all Data types

Fast new data sources onboarding

Traditional and Advanced Analytics use cases

**Modern Data platform**

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Challenges – Questions Need to be Answered

- **What workloads do we need to support?**
  - Batch vs Streaming
  - AI vs Traditional analytics
- **What protocols need to be supported?**
  - HDFS vs S3 vs …
- **Where best to run your data platform?**
  - On-premises vs Cloud vs Hybrid
- **Data considerations**
  - Gravity
  - Sovereignty
  - Compliance
  - Security

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# The Skills Challenge

- **Bringing in new technologies**
  - Selecting the right tools for the right workload out of the huge number of choices
  - Containerization

- **How to get support for new tools and technologies?**
  - Fast moving ecosystem
  - Many open source projects

- **Where do we find the people for these platforms and workloads?**
  - In demand skills
  - Upskilling existing teams

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Building for the Future

- Can we futureproof your data platform?

- How to not make the same mistakes again?

- Does everything need to move to the cloud?

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Modernization Challenges

## Data sprawl

- Data is spread across multiple on premises and public cloud locations

- Data is accessible via multiple protocols (NFS, HDFS, S3)

- Finding relevant data

- Managing multiple copies of data

## Data governance and data gravity

- Data classification

- Data sovereignty

- Regulatory compliance

- Not all data can move to public cloud--leverage data catalog to ensure compliant data movement and data placement

- Expensive lift and shift

## Performance, scalability, and durability

- Bring data closer to compute -- long latencies when accessing data from data lake storage

- Cost prohibitive to keep all data in high performance storage tier

- With non-persistent cache, all data must be reloaded in the event of failure

- Elasticity and cloud bursting
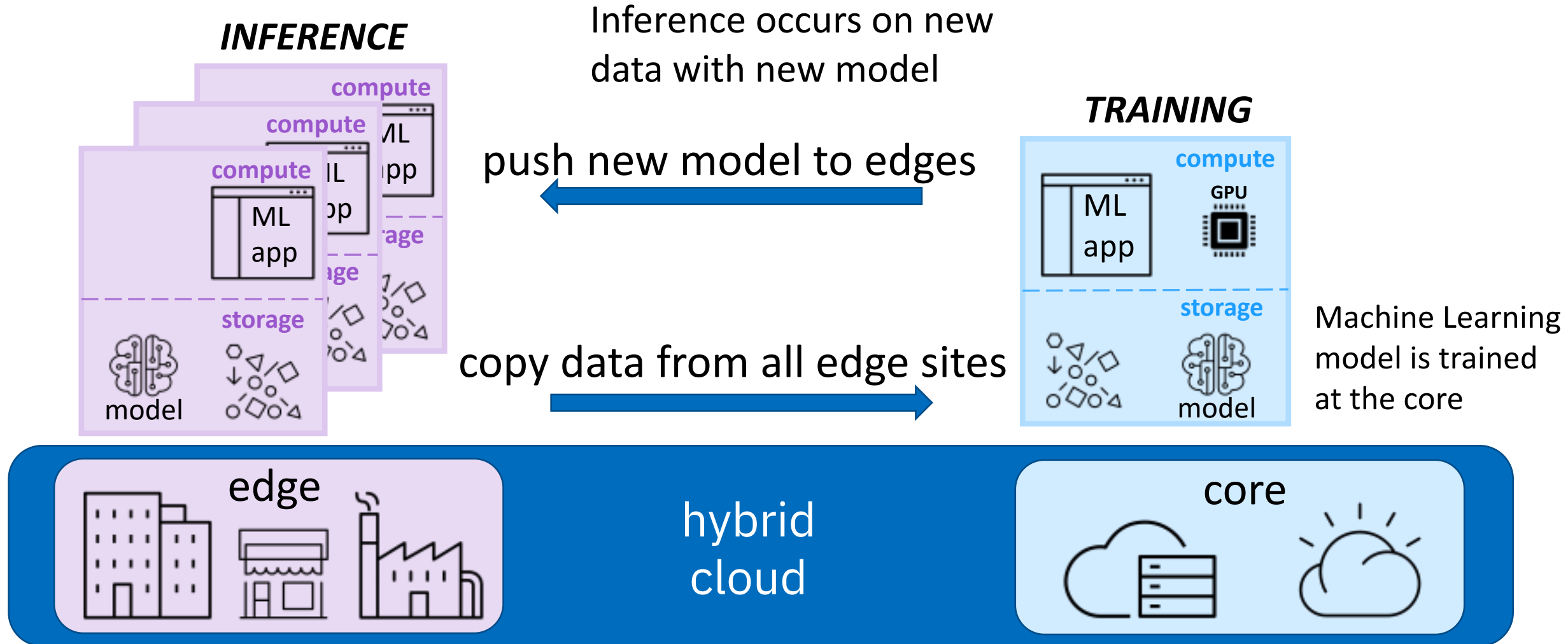
- Computational storage

## Data security

- Encryption of data in flight and at rest

- Hybrid key management

- Role based access control

SNIA CSTi | CLOUD STORAGE TECHNOLOGIES

# Evolving Workloads

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# Machine Learning with Data from All Edge Sites

What usually happens today…



**INFERENCE**

Inference occurs on new data with new model

push new model to edges

**TRAINING**

Machine Learning model is trained at the core

copy data from all edge sites

edge

hybrid cloud

core

SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Federated Learning

Models are retrained on the edge

**TRAINING and INFERENCE**

**compute**
**compute**
**compute**

ML app
ML app
ML app

**storage**
**storage**

model

Inference occurs on new data with new model

push new model to edges

request retrain

pull models from edge

(not raw data)

Machine Learning model is trained on the core

**TRAINING**

**compute**

ML app

GPU

**storage**

model

Aggregator collects models from each edge site, retrains, and redistributes new models

edge

hybrid cloud

core

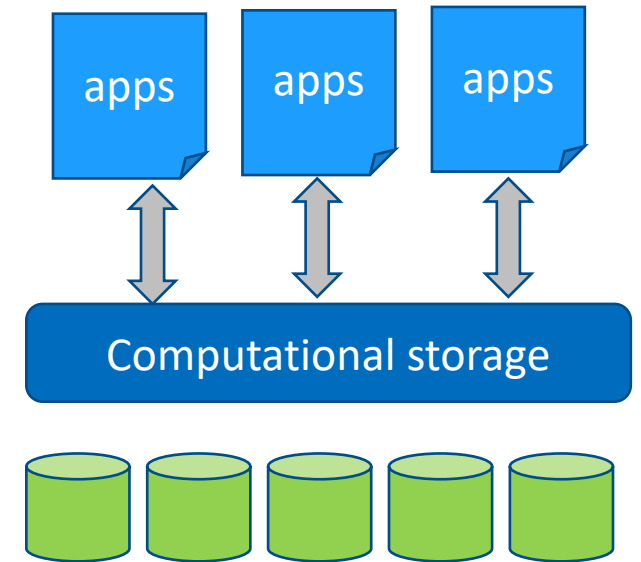SNIA. CSTI | CLOUD STORAGE TECHNOLOGIES

# Value of Federated Learning

- Improve model training across locations
- Address data privacy, locality and security
- Adhere to regulatory compliance
- Tackle data volumes at lower cost and risk (e.g., minimize egress charges)

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# Looking Towards the Future

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# Look Towards the Future – Storage for the New Big Data

- True hybrid cloud data fabric

- Acceleration technology: FPGA, GPU, DPU, IPU,…

- Computational storage

SNIA CSTI | CLOUD STORAGE TECHNOLOGIES

# Looking Towards the Future

- Data and Analytics as a Core Business Function

- Data and Analytics at the Edge

- Operationalization of AI
  - DevOps, AI Ops, ML Ops

- The Data Lakehouse
  - Bringing together the best of the data warehouse and data lakes

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES

# Thanks for Viewing this Webcast

Please rate the webcast and provide us with feedback

This webcast and a copy of the slides will be available at the SNIA Educational Library https://www.snia.org/educational-library

A Q&A from this webcast will be posted to the SNIA Cloud blog: www.sniacloud.com/

Follow us on Twitter @SNIACloud

SNIA
CSTI | CLOUD STORAGE TECHNOLOGIES

# Thank You

Questions?

SNIA. | CLOUD STORAGE
CSTI | TECHNOLOGIES