

Regional SDC Denver April 30, 2025

Cloud Storage Considerations for Retrieval Augmented Generation (RAG) in Al Applications

Daniel Falkner - Senior Product Manager - Azure Object Storage Scott Hoag - Principal Product Manager - Azure Object Storage



Agenda







Object Storage: Workloads over the years





Al Workloads & Storage Requirements



Azure Storage portfolio

Durable, highly available, massively scalable

Block storage

Services

Azure Disk Storage Azure Elastic SAN Azure Container Storage Unique capabilities

Ultra Disk and Premium SSD v2 Azure Elastic SAN Shared disks

Object storage

Services

Azure Blob Data Lake Storage

Unique capabilities

Premium Blob Multi-protocol access (e.g. NFS, HDFS) Azure Managed Lustre

File storage

Services

Azure Files Azure NetApp Files

Unique capabilities

Native NetApp File Storage Azure File Sync

Capacity

100s of trillions of objects across many exabytes of data

Throughput

>500 Tbps average
(>150 exabytes per month)

IOPS
>500M tps
(>1 quadrillion per month)



Azure Blob Storage: Multi-protocol, single storage platform





Al Pipeline – Storage centric view





Al Fine-Tuning – Storage centric view





Al Pipeline - Storage Requirements



Requirements

Training / Fine-Tuning

- Ingestion: Bring raw training data to Azure
- Data Preparation: Integration with Spark, MosaicML, etc.
- **Training/Fine-Tuning:** Data to GPU nodes, checkpoints to storage. Integration with PyTorch and other ML frameworks
- Data Management: Secure & cost-efficient retention

Deployment/Inference

- Deployment: Model distribution and load times
- **Data Management:** Model versioning, retention of inference inputs and outputs



9 | ©2024 SNIA. All Rights Reserved. © 2024 Microsoft Corporation. All Rights Reserved

RAG with Object Storage



Bringing domain knowledge to LLMs





RAG: <u>Retrieval-Augmented</u> <u>Generation</u>

Combine reasoning + knowledge



RAG Pipeline – Storage centric view





RAG Pipeline – Storage Requirements



Microsoft

RAG with Blob Storage





SSD-backed Premium Blob Storage for RAG

	Latency (400KiB PDFs)
Operation	Premium Blob (SSD) vs Standard (Hot) Blob (HDD)
PutBlob (Data Ingest + Chunk Writes)	Premium ~2x faster
GetBlob (Vectorization)	Premium ~3x faster
GetBlob (Retrieval)	Premium ~3x faster

Premium delivers ~3x faster RAG performance with 65% savings on Transactions!

16 | ©2024 SNIA. All Rights Reserved. © 2024 Microsoft Corporation. All Rights Reserved



Blob Storage: Scaled accounts



Key Takeaways

Leveraging Object Storage for building AI Apps ...









Resources

- <u>Building AI applications that leverage your data in object storage | BRK216</u>
- Open AI accelerates AI innovation with Microsoft Azure Blob Storage scaled accounts
- Training and fine-tuning your large AI models using Azure Blob Storage
- Bring your data to AI applications RAG with Azure Blob Storage
- SNIA SDC 2024 Supercharging OpenAl Training with Azure Blob Storage
- OpenAI transforms AI model development with Azure Blob Storage | Microsoft Customer Stories