

Object Storage Key Role in Future of Big Data

Anil Vasudeva

President & Chief Analyst





- The material contained in this tutorial is copyrighted by the SNIA and the author.
- 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
 - SNIA and author 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 Education Committee.
- 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.





Object Storage: Key Role in Future of Big Data

This session will appeal to Data Center Managers, Development Managers, and those that are seeking a fundamental understanding of evolution of Object Storage to deal with the explosive growth in unstructured data. The session delves into the impact of Object Storage architectures on next generation storage in Private and Hybrid Clouds as well as its ability to gain meaningful business insights from Big Data given Object Storage flat hierarchical structure and associated metadata attributes governing policies applicable to Big Data. The metadata in object storage can become key in application driven Software Defined storage. The rapid adoption of Object storage by both startups and major IT vendors is a testament to their advantage as the next generation storage for Big Data. The session is targeted to bring a clear understanding of this key technology to developers as well as systems admins and large service providers providing Big Data Analytics.

Unstructured Data Explosion





Object Storage Key Role in Future Data

Harnessing Big Data Infrastructure





Big Data Analytics



Big Data Ecosystem

System Tools



Generation	Operational IT	Analytics	Usage
Data Class Types	Store Access Prepare	Analyze Visualize	Analyze Business
Data Types - Structured (Relational) - Unstructured (Adhoc)	Data Mgmt & Storage - Store - Secure - Access - Network	Data Analytics - Algorithmics - Automation - In Real Time	Business Analysis - Decision Support - Just InTime Business Model
Data Class - Human - Machine Data Velocity	Engines - Hadoop/MapReduce - Apache Tools - Cloudera/IBM/EMC - Visualization	Business Analytics - Visualization - Interoperate with SQL- RDBMS - BI/EDW	Business Use - Market Penetration Enhancements Cash Flow/POI
- Batch - Streaming	Prepare Data For Analytics - ETIL / Data Integration - Workflow Scheduler		- Cash Flow/ROF

Big Data HDFS Architecture





HDFS - Actively Maintaining High

HDFS

- Immutable File System Read, Write, Sync/Flush No random writes
- Storage Server used for Computation Move Computation to Data
- Fault Tolerant & Easy Management Built In Redundancy, Tolerates Disk & Node Failure, Auto-Managing addition/removal of nodes, One operator/8K nodes
- Not a SAN but high bandwidth network access to data via Ethernet
- Used typically to Solve problems not feasible with traditional systems: Large Storage Capacity >100PB raw, Large IO/computational BW >4K node/cluster, scale by adding commodity HW, Cost ~\$1.5/GB incl. MR cluster

Object Storage Key Role in Future Data

Big Data - Market Requirements



Types of Data Organizations Analyze 59% Customer/member data 68% 44% Transactional data from applications 68% 69% Application Logs 37% 64% Other Types of Event Data 23% 41% Network Monitoring/Network Traffic 33% **Online Retail Transactions** 51% **Other Log Files** 26% 28% Call Data Records 46% Web Logs 21% 36% Text data from social media and online 15% 36% Search logs 11% 18% Trade/quote data 15% Hadoop 18% Intelligence/defense data Non Hadoop 11% 21% Multimedia (audio/video/images) 9% 8% Weather 3% Smartmeter data 6% 3% Other (please specify)

Object Storage Key Role in Future Data

Big Data Savings with Open Source



Legacy BI vs. Open Source Big Data Analytics



Targets for Object Storage





Instantly available On Boot Ups Most Accessed Videos Very Read Intensive Rendered (Texture & Polygons) Stored Data-Very Read Intensive Wall St/Financials/ Federal/State/CityData Geo-Resources Maps

Object Storage - Goals





Object Storage Goals



Reduce cost and complexity

of unstructured data storage infrastructure.

Auto-manage unstructured data

documents, emails, presentations, audio files, messages..

• WORM

Upload once, access from any geographic region

• Unlimited File Size.

Upload files of virtually any size

Data Consistency:

Changes made to a file in one location immediately available at all locations

• True Global Namespace:

Leverage auto geo-positional routing for increased performance and consistency.

• Self-Healing Data Integrity

Continuous hash check, if problem detected, system repair the file automatically.

• Self-Provisioning:

Provision your own storage as needed.





• Object-based storage supports shared tenancy.

Objects have their own custom metadata functioning as fairly autonomous data instances - can carry access policies to support multi-tenancy - a big aspect of the relevance of object storage to public cloud infrastructures.

• Object storage similar to low-cost, grid like architectures.

Object storage is well-suited to run on a loosely federated set of low-cost, industry-standard servers acting as nodes in a cluster. Cluster easy to expand by adding nodes - automated data layout rebalance the system to incorporate them. This drives the low acquisition cost and lower OPEX..

Why Object Storage



Massive scalability

Within specified policies system can scale nearly indefinitely in number of files or objects and capacity of system..

• Large metadata > custom control over data.

Object Storage's expanded Metadata automatically and consistently tackles traditional storage challenges such as tiering, security, migration, redundancy, and deletion vs File Storage's only file type, creation date, and last-accessed date.

Data is inherently Immutable

Saved object data is tagged with a unique ID, guaranteeing the immutability of that object. To modify an object, you create a new, and keep versioning history - a natural fit for archiving and records management guaranteeing that data has not been modified or tampered with.

Why Object Storage for Big Data



Object Storage Characteristic

Object is API Driven Storage

- Faster Development,
- 10x reduction in Code
- REST APIs speed Development
- Location Transparency
 - One Storage System/Access Point across
 - Multi-global apps
- Self Managed Storage - No LNS, Never recode when syst Change
- No limit on number of Objects
- Object size can be up to 5 TB
- Central Data storage All systems
- High Bandwidth
- > Five 9s of Durability
- Versioning Lifecycle policies
- Integration with industry standards
- Typical Suppliers:

Open Source/Swift, AWS/S3, Glacier, EMR,Redshift; Basho, Caringo, Cleversafe, Cloudian, Coho Data, DDN, EMC, Exablox, HDS, Huawei, IBM, NEC, NetApp, RedHat, Scality, Storiant, SwiftStack, Tarmin



Linear Scalability Scales to billions of objects



Support for large files Object sizes are in TBs



Web friendly Firewall friendly, http, REST accessibility

metadato

Objects can be extended

Metadata and

to multiple policies (Immutability, retention,

extensibility





Geo-scale Geo-replicated and distributed

Object Storage – Global Reach



GLOBAL REACH & DATA LOCALITY



HYPERSCALE PERFORMANCE & CAPACITY

Replicate & Collaborate (Ingest, Access & Update Data at local speeds across multiple locations

Object Storage Key Role in Future Data

Data Integrity of Big Data Storage



Issues

Preventing Data Loss Maintaining 'Always On' System Scaling Storage Capacity while Controlling Costs Preventing Data from unauthorized Access



Solutions

Choose Object Based storage vs. File Based Systems To manage billions of objects Information Dispersal Algorithm Using Error Codes Preventing Data from unauthorized Access Scaling Storage Capacity while Controlling Costs





Protocols

Data Storage Tiers Evolution





Object Storage Key Role in Future Data

File vs Object Storage Characteristics SNIA

2012-2016 Workload Growth

Enterprise IT must adapt to support mixed workloads



- Optimized for structured/semi-structured
- CRM, ERP, Database apps, email, etc.
- Built for performance (10s of ms latency)
- Directory structure, location-dependent
- Manual IT administration & provisioning
- Writes to file require exclusive lock
- Limit on number of files in a directory
- No user meta-data
- Large files hard to seek
- File create require directory exclusivity lock
- Not Web or firewall-friendly

DodeJS 2016 48 Million 2012 6 Million 700% 0bject/Cloud Storage

Next-Gen Web, Mobile, Cloud Applications

- Optimized for unstructured
- Archiving, Web, mobile, cloud apps
- Built for scale/efficiency (100s of ms latency)
- Location-independent Metadata & objects stored together
- Multi-tenancy self-service access
- · Object supports multiple writes, no locking
- Objects are limitless in size, 1 MB to TBs
- Objects support extensible meta-data
- · Objects can be viewed with no limitation
- No locking required to create files
- · HTTP, REST-based access, Web & firewall-friendly

Object Storage Key Role in Future Data

Approved SNIA Tutorial $\ensuremath{\mathbb{C}}$ 2015 Storage Networking Industry Association. All Rights Reserved.

Hierarchical vs Flat Storage Schema



SN

Global Education

Metadata: Key to Object Storage Adoption SNIA

VS.



File Name: CATSCANJQSMITH Created By: Technician 1 Created On: 01-01-2001 File Type: .DICOM

Object



Object ID: 12345 File Type: .DICOM Patient Name: John Q. Smith Patient ID: 555-55-5555 Procedure Date: 01—1-2001 Physician Name: Dr. Organ Physician Notes: .WAV File Prior 1: XYZ.DICOM Modality: XYZ Manufacturer: XYZ Diagnosis: XYZ Description: XYZ Custom Metadata: XYZ

Object Storage Key Role in Future Data

Object Storage Positioning







- Object Storage:
 - Data repository of billions of Objects
- Key Architecture:
 - Šeparation of Metadata and Data

• Features :

- Flat Namespace
- Infinite Scalability
- Elasticity
- Cost-efficiency
- Data durability
- Distributed System
- No Single point of failure

App Driven SDS = Metadata Driven Object Storage SNIA





In implementing object storage:

- A unique identifier assigned to each piece of data
- Identifier allows servers or users to get access to Object without any knowledge of where data is physically stored.
- Data can be broken into chunks, stored anywhere worldwide, can move within a system and have multiple copies.
- Object tag ensures that right piece of data is returned and verifies integrity of that data.

Open Stack Object Storage





Object Storage Key Role in Future Data

Centralized Storage using Public/Private Cloud





Object Storage: Key for Big Data



Object Storage Characteristic

- Object is API Driven Storage
 - Faster Development, REST APIs speeds Dev.,
 - 10x reduction in Code.
- Location Transparency
 - Access Point across Multi-global apps
- Self Managed Storage - No LNS, Never recode when syst Change
- No limit on number of Objects
- Object size can be upto 5 TB
- Central Data storage All systems
- **High Bandwidth**
- > Five 9s of Durability
- Versioning Lifecycle policies
- Integration with industry standards



Linear Scalabi Scales to billions objects

Web friendly

REST accessibility



Support for large files Object sizes are in TBs



No Locking No lock on write or create operations



Metadata and extensibility Objects can be extended to multiple policies (Immutability, retention,



Geo-scale Geo-replicated and distributed

Typical Suppliers:

Open Source (Swift), AWS (S3, Glacier, EMR, Redshift) Basho, Caringo, Cleversafe, Cloudian, Coho Data, DDN, EMC, Exablox, HDS, Huawei, IBM, NEC, NetApp, RedHat, Scality, Storiant, SwiftStack, Tarmin



- Object Storage is rising to meet the needs of next generation storage bringing in infinite scalability, ease of management and security in Cloud Computing and opening greater insights into Big Data by leveraging extensive attributes in Metadata associated with Object Storage while significantly lowering cost of storage using volume x86 hardware to manage both servers and storage lowering CapEx and OpEx for enterprises and cloud service providers.
- Several start-ups and major IT vendors are now endorsing it.



The SNIA Education Committee thanks the following Individuals for their contributions to this Tutorial.

Authorship History

Original Author : Anil Vasudeva, IMEX Research.com

Additional Contributors

Please send any questions or comments regarding this SNIA Tutorial to <u>tracktutorials@snia.org</u>