



I D C T E C H N O L O G Y S P O T L I G H T

Object-Based Storage: The Need for Decentralized Scale-Out Architectures

February 2015

Adapted from *IDC MarketScape: Worldwide Object-Based Storage 2014 Vendor Assessment*, by Ashish Nadkarni, IDC #253055

Sponsored by NetApp

The growth in unstructured data continues to push the limitations of traditional storage architectures. This has led to newer highly scalable technologies such as object-based storage (OBS) gaining momentum. Object-based storage (OBS) solutions offer the ability to access and manage data using RESTful protocols that enable scaling via a geographically dispersed namespace. Unlike traditional scale-up architectures, which largely rely on the capabilities of the controllers to dictate scalability in terms of performance and capacity, OBS enables organizations to add processing power or storage capacity independently, as applications may require. The ability to scale capacity across multiple datacenters, while maintaining a single virtual view, also ensures accessibility and resiliency. This Technology Spotlight examines the reasons behind the growth in adoption of OBS architectures, and also looks at the role that NetApp's StorageGRID Webscale plays in this strategically important market.

Introduction

Simply put, object storage is a type of storage architecture in which data is stored as objects in logical containers and not in a file system hierarchy or as blocks of data on hard drive sectors or magnetic tape tracks. It's a higher level of service, and allows the data to be associated with rich metadata. The object contains not only the original data, but also an additional amount of metadata that describes the data, a digital fingerprint (which enables verification that data has not been altered in any way and ensures correct data is given back to the application), and a global unique identifier, which makes the data searchable no matter where it's stored. (A global unique identifier enables applications to access data regardless of physical location in a truly global namespace.)

Metadata describes the data in the object (metadata may be stored with the object or separately). Metadata can be used for a number of things, including searches, but more importantly it can be used to trigger storage policies that store objects in a particular location, storage tier, etc. If an object is changed, the global unique identifier also changes, creating a new object. The original object with its data remains retrievable and searchable (in such OBS solutions, objects can be made immutable and a change to an object forces the creation of a new object, thus ensuring the original object remains untouched by the change).

Because OBS uses the language of the Web — Hypertext Transfer Protocol (HTTP) — it is capable of describing the location of stored data anywhere in the world, making it accessible from anywhere in the world. Such platforms offer a unique value proposition in that the data persistence and access is independent of the application. In other words, a stateless connection to the storage allows the creation of cloud-like applications.

Early object-based platforms suffered from "necessity crisis," and were too cumbersome to deploy, in some cases causing platform lock-in because of their proprietary access mechanisms. In spite of

their from-the-ground-up design, and their departure from how traditional SAN and NAS arrays are deployed and a lack of standard interfaces, object-based storage had been difficult to deploy. Thanks to the growing popularity of protocols such as Amazon's S3, OpenStack SWIFT, and SNIA CDMI becoming de facto RESTful access interfaces, this situation has changed. This has, in turn, spurred further investments in making OBS solutions simpler to deploy and ubiquitous.

By comparison, increasing capacity in traditional storage infrastructures means setting up additional servers, hard drives, and tape backups, along with increasing personnel to manage it all. Storage administrators have also been burdened by estimating how much storage to provision. Getting that capacity amount wrong has meant not having the needed storage or overspending on unnecessary capacity that then sits idle.

With OBS, resiliency can be built-in, as one object can be automatically replicated and stored on multiple devices across datacenters. For data that is accessed infrequently, life-cycle policies can be configured to active and archived objects.

Object access control is also dictated by identity and access control lists. Uploads and downloads are safeguarded with secure-socket layer (SSL) endpoints, as well as encryption for data at rest.

The Need for Object-Based Storage

New deployment and delivery models such as online or cloud-based content delivery, mobile workloads, digital archiving, and Big Data and business analytics are driving a need for scale-out architectures.

The sheer amount of data being created by digital systems is growing in almost incomprehensible proportions. In the past two years, more data has been created than in all previous history. Today, more than 4.4 zettabytes of data exist, but by 2020, that will grow tenfold to 44 zettabytes, according to the *IDC Digital Universe Study*.

Over the next five years, businesses will more often be moving toward petabyte-scale data storage, and by way of their design, OBS solutions are turning out to be the right choice for balancing scale, complexity, durability, and costs for enormous data sets.

The continued expansion of business-critical information, diverse data sets, new mobile and cloud service platforms, and rich content within extended enterprises continues to change the storage dynamic in a wide range of industries and organizations. Regulated industries, such as healthcare and financial services, are often challenged to store data sets for extended periods while still being able to quickly retrieve them.

For example, in financial services and healthcare, documents and images are required to be kept for years — even decades. A mammogram must be kept a minimum of five years, but, depending on the state, may be required to be stored for the life of the patient. Under the Securities Exchange Act, Rule 17a-4, broker-dealers must keep securities transaction records accessible at least six years. When a regulatory/legal challenge arises for a bank or brokerage, or a patient's medical history is required, the data must be retrievable with little notice. This is especially the case for older data sets — i.e., the durability and integrity of the objects are critical. Offline mechanisms like tape cannot offer this level of integrity assurance.

In the retail space, organizations collect data on customers, which can be used to analyze general populace trends or individual habits. As the amount of data accumulates, it's often lost in time in irretrievable archives. Even data that retailers are unsure of how to use today may be useful as new analytics platforms arise in the future, giving enterprises good reason to store too much data instead of losing what could someday prove a goldmine of information.

Benefits Unique to Object-Based Storage

Because OBS uses a single, flat namespace for an object, replacing both path and filenames with object identifiers, object addresses are nearly infinitely expandable. The object global namespace transparency also makes it quick and easy to locate for users and distributed applications. Because applications can be programmed to associate metadata with the stored objects, there is no need for a database on the application side.

With OBS, files become immutable objects that cannot be manipulated because each is created with its own unique digital fingerprint, ensuring its integrity. If someone attempts to overwrite it, an OBS system will create a new object with its own unique identifier, keeping the original file intact.

With good OBS systems, an object's digital fingerprint of data is checked every time it's retrieved. OBS solutions like these are also beginning to offer an additional level of protection: If a user tries to retrieve a file that has become corrupted, the system doesn't just offer an error message, but it also recovers an accurate object from backup while also repairing the bad one.

Additionally, OBS systems can be set up to perform data health checks periodically to determine if any files have become corrupted, while at the same time repairing them through error-correction coding.

As data storage paradigms shift toward private, public, and hybrid clouds, OBS enables IT organizations to shift IT resources from maintenance to new initiatives. Simply put, OBS deployments offer flexibility to both grow and maintain storage, while also offering 24/7 access to data from any location and to any computer platform, whether it's a traditional workstation or a mobile device.

Objects are simple, easily distributed and replicated, and can be quickly and globally accessible in the cloud. That makes objects perfect for active global archives and accessible through mobile applications.

Most enterprise hardware being shipped by vendors today is in support of unstructured data. Unitary file servers and scale-up appliances and gateways continue to decline due to capacity limits, while scale-out or clustered file solutions continue to grow.

At the same time, a generation of Web 2.0 vendors is creating new applications designed for use with object-storage's RESTful interfaces, mainly Amazon's S3 and OpenStack SWIFT.

Additionally, industry support from the Object Storage Alliance and SNIA, with nearly two-dozen storage vendor members, is helping educate enterprises on the use cases and the problems solved by OBS.

Key Characteristics of OBS Platforms

An object platform is no good if it does not have adequate intelligence to handle "objects" of different types. In other words, it needs to have extensive capabilities to extract the "DNA" from the object, which allows it to understand what this object is, place it on the appropriate storage tier, and maintain the proper level of metadata detail.

One of the key differentiators between true OBS platforms and scale-out file-based solutions with object interfaces is how they manage data resiliency and dispersal. The more robust the object platform in terms of its dispersal capabilities, the better it is suited for cloud and cloud-like deployments in which data may be ingested and/or analyzed concurrently across multiple locations globally and dispersed via replication and/or erasure coding on a granular basis (per object or container). In addition, the platform needs to have robust clustering, data persistence, replication, and

conflict resolution capabilities. Earlier variants of OBS lacked this capability and therefore could not function in a multisite configuration.

Look for platform scalability — not just hardware, but throughput, file size, and from file volume perspectives. An OBS system that's appropriate for a given environment will allow each dimension to scale independently.

As with any storage system, data layout and organization is crucial for mining existing data for patterns that may build new business cases. Users should consider a solution that supports advanced metadata, indexing, and analytics.

As storage architectures continue to grow exponentially in size, there is a greater need for efficient data management and reduction techniques. Because object storage uses metadata to describe a file, that information can be used to create automated data management policies. Key policies a platform should offer include data protection through striping, compression, deduplication, and scheduled deletion.

Data optimization technologies that include automated data tiering will also be key. An OBS solution appropriate for a given environment will allow many, if not all, of the previously mentioned features.

It's also important to determine a vendor's commitment to the platform now and for the future, and to consider the vendor's past track record for incorporating new features over time. Established vendors have greater R&D resources to develop new feature sets as a platform matures.

Look for a vendor partner ecosystem, too, for applications and on-ramping. Unlike SAN and NAS, OBS has a less-established set of interfaces for existing enterprise applications. Therefore, the more comprehensive the ecosystem, the better placed the supplier is to offer an end-to-end workload-optimized or use-case-focused solution.

Considering NetApp StorageGRID Webscale

NetApp has been selling StorageGRID Webscale as a software-defined OBS platform targeted at large-scale distributed content repository deployments. While StorageGRID is available as a software-only solution, the most common deployments are with NetApp E-Series storage, where StorageGRID effectively utilizes E-Series DDP (dynamic disk pools) for node-level erasure coding, performance, and density. NetApp added the StorageGRID Webscale appliance for easier, more dense deployments. The appliance and software-only version can be mixed as needed in a single deployment. The appliance comes in a 4U/60 drive form factor.

NetApp StorageGRID Webscale is designed to eliminate the constraints of mapping data into predefined containers as blocks or volumes, allowing repositories to accommodate billions of files or objects and many petabytes of capacity in a single, unbounded repository that can span platforms, technologies, and sites. Webscale offers geo-distributed Erasure Coding controlled by a dynamic policy engine, enabling objects to move between storage tiers, locations and to be replicated.

StorageGRID Webscale 10.1 has added several new features, including support for Amazon S3 as a data placement storage tier. It also supports a single policy engine that controls data placement across private and public clouds.

By providing a metadata-based policy management approach, StorageGRID Webscale transcends the limitations of traditional storage containers and hierarchical organization structures, allowing data to be organized, accessed, replicated, and managed using multiple user-defined and metadata-driven criteria.

Data placement policies can be applied years after ingest, allowing business owners to adapt their protection strategies over time as business needs/rules change.

NetApp StorageGRID Webscale is suitable for single-datacenter or multi-datacenter deployments with many sites across the globe. Key differentiators of the StorageGRID Webscale platform are native support for S3, CDMI, and tape as well as the ability to continuously evaluate policies for an object at rest. The StorageGRID Webscale platform is strong in its approach to multisite configurations, a powerful policy-based management framework; security features such as tamper detection, robust auditing, and its streaming architecture (versus a store and forward architecture) that enables reduced time to first byte; and efficient reads even when data is compressed or encrypted.

The StorageGRID Webscale policy engine can drive a lower total cost of ownership. Unlike other policy engines in the market, all application metadata (including S3 metadata) is available to the NetApp policy engine. Policies can be applied not only at ingest time, but at any other time, and new policies can be defined for objects at rest. This provides the flexibility to monitor and meet cost and performance SLAs dynamically, instead of being locked into a particular policy. It allows enterprises to meet and adapt to continuous changes in their business environments.

IDC placed NetApp in the Major Players category in its *IDC MarketScape: Worldwide Object-Based Storage 2014 Vendor Assessment* report. This position reflects NetApp's significant investment and product deliveries over the past year, as well as its quest to further penetrate this market and gain traction among both its existing customer base and its newer customers by leveraging its existing position in the storage industry.

Challenges

Historically NetApp's StorageGRID hasn't had the same level of visibility inside of NetApp compared with its other storage platforms and solutions. This changed dramatically last year with the increased investment in object storage and the subsequent release of StorageGRID Webscale, as well as the addition of a field-facing organization focused on delivering NetApp's newest solutions (Emerging Products Group).

To take advantage of the opportunity provided by the launch of its newest OBS, NetApp needs to find specific use cases for StorageGRID Webscale that cannot be satisfied just by deploying its Fabric-Attached Storage (FAS) or E-Series arrays.

The company needs to continue to convince the market that its OBS platform is not a fulfillment mechanism for its other storage offerings, and that it's indeed serious about this emerging software-based storage market segment. NetApp also has to go above and beyond to convince the market that it offers the same level of support that its customers enjoy for its principal platforms. The fact that NetApp has added AutoSupport to its OBS is a step in the right direction — to bring many of the enterprise-grade features to OBS solutions.

The value of the object platform lies in the number and diversity of solutions it enables. It can be solutions for unstructured data, geographically dispersed data and computing, and semistructured data, as well as across industries like healthcare, energy, manufacturing, and research to name a few.

Given NetApp's success in the storage market, it is well-positioned to continue to further improve its reputation in the OBS market in the near future.

Conclusion

Outside the datacenter, storage traffic is undergoing a dramatic shift as the growth of mobile, social, and cloud is heavily tilted toward the use of IP-based connectivity mechanisms for consuming storage resources. As businesses move toward petabyte-scale data storage, object-based storage architectures are turning out to be the right choice for balancing scale, complexity, and costs.

ABOUT THIS PUBLICATION

This publication was produced by IDC Custom Solutions. The opinion, analysis, and research results presented herein are drawn from more detailed research and analysis independently conducted and published by IDC, unless specific vendor sponsorship is noted. IDC Custom Solutions makes IDC content available in a wide range of formats for distribution by various companies. A license to distribute IDC content does not imply endorsement of or opinion about the licensee.

COPYRIGHT AND RESTRICTIONS

Any IDC information or reference to IDC that is to be used in advertising, press releases, or promotional materials requires prior written approval from IDC. For permission requests contact the Custom Solutions information line at 508-988-7610 or gms@idc.com. Translation and/or localization of this document require an additional license from IDC.

For more information on IDC visit www.idc.com. For more information on IDC Custom Solutions visit http://www.idc.com/prodserv/custom_solutions/index.jsp.

Global Headquarters: 5 Speen Street Framingham, MA 01701 USA P.508.872.8200 F.508.935.4015 www.idc.com