# File Systems and the Cloud for Digital Archival of Motion Picture Assets

## 1. Introduction

Motion picture productions generate vast amounts of digital data. To provide a sense of scale, imagine a motion picture camera capturing at 24 frames per sec. In the span of one minute, it generates 1440 frames, resulting in about 130,000 images for a 90-minute feature. Now, include the data from multiple cameras, multiple takes, and include audio, editorial, visual effects, multi-language audio tracks and subtitles, plus ancillary data and it is common to have many millions of files generated during the life cycle of a motion picture. In current day 4K productions, it is commonplace to see this dataset consume multiple petabytes of storage.

While some of the data is considered temporary and is not typically archived, digital storage systems are still used to read, write and access this data. Choosing what data to archive is a choice that content owners have to make, in consultation with technologists and archivists, and is dependent on numerous factors such as cost, retrievability, reliability, security and utilization. For a discussion on selecting assets for digital preservation, please refer to this talk with archivists (https://academydigitalpreservationforum.org/2021/11/19/archivists-talk/).

In a separate paper, we examined the access, retention and security considerations for any digital archival storage, and looked at how tape and hard disk drive satisfy these considerations.

In this paper, data storage is explored from the standpoint of an enterprise and how data is accessed at scale, going from individual files to large distributed file systems and building up to the cloud. Additionally, the capability of value-added scenarios offered by the public cloud providers is explored.

## 2. Data Considerations

In an earlier paper, *Digital Storage Considerations and Devices for Archiving Motion Picture Assets*, we reviewed the different requirements for our data... For reference these requirements are briefly duplicated below:

### 2.1. Data Retention

The digital data needs to be preserved for an indefinite period of time. The fidelity of the data needs to be maintained without any compromise, irrespective of whether it is short-lived or long-lived. It is important to note that the data needs to be preserved even in the face of catastrophes such as natural disasters and wars.

### 2.2. Data Access

The digital data needs to be accessible and data access requirements vary across different types of assets. These requirements can vary from low-latency to high-latency. Additionally, data has to be accessible over time, and this means that the data is stored on media that is supported and operable as technology changes. For the stored data to be accessible, it then becomes important for the media to be backward compatible, or it is necessary for media to be migrated to the new media or even the new data format.

## 2.3. Data Protection and Security

The digital motion picture data needs to be protected. The data needs to be protected from corruption and failure of the storage media, but also from any disaster that might strike the storage location – rendering access to the storage impossible, and thus loss of the archive.

The data needs to be secure, either from unauthorized access by bad actors, or just improper or erroneous access by authorized users. In many cases, it may be important for the data to be encrypted, so that if some chance access is obtained to the storage media, an additional layer of protection is available before the data is stolen.

## 2.4. Data Utilization

The storage lifecycle of content for motion pictures typically starts with data that is minimal (at concept stage), grows to very large (at capture stage) and systematically reduces in a deliberate process (during editorial) to create a highly specific viewing experience (VFX and color correction) and then rapidly grows again to support all the distribution channels and multiple languages. From this it is evident that depending on the stage of the motion picture, while data is important to retain, not all the data is archival, and not all of the datasets are equally important in an archival context.

Data utilization should be viewed in the context of what stage a production is, and the access, security, protection and storage needs for that stage. The storage architecture for that stage should then to be constructed in a fashion that optimally supports those needs.

## 2.5. Other Considerations

In addition to the basic data considerations that are listed above, there are numerous enterprise-level requirements that need to be satisfied for any digital data archive. Some of these are:

1. Common ontologies and metadata in formats that are flexible, extensible and robust are important to support discoverability and automation. Similarly, file formats of files destined for archive should ideally not be in proprietary formats.

   It is important that the production dataset's ontology is maintained, so that entity relationships within the dataset are available and accessible to anyone who accesses the digital archive. This ontology must be flexible, extensible (as per different production needs), robust and able to connect and map to other evolving and extending ontologies.

The metadata needs for a motion picture are extensive and vary drastically based on the production group that is working on it. This metadata needs to be filtered appropriately for the group that needs it, yet at the same time needs to be flexible, extensible, robust and easily transportable.

2. Provide a secure mechanism for storing data so only authorized applications can access assets via a common search and discovery interface.

3. Asset management driven by enforceable policies regardless of the infrastructure on which they are hosted or the applications that manage them.

4. And many other considerations that are related to
   a) Business Processes and workflow
   b) Corporate Governance
   c) Clearance, Rights and Legal issues relating to ownership
   d) Etc.

## 3. File Storage – Local and Distributed

In *Digital Storage Considerations and Devices for Archiving Motion Picture Assets,* the underlying technologies and mechanisms of storage for tape, hard disk drives (HDD) and solid-state drives (SSD) were explored. From the device (magnetic level or integrated circuit level) we now move up to the data level. This is the level at which the binary bits (0s and 1s) are consolidated into the digital file we are familiar with. These files need to be organized, secured, and have cross relationships with other files in the production chain. This is where operating systems and file systems come in, enabling a translation from the structure that is best suitable for a user to the actual storage on the device.

While storage devices store data, it is in a format that is optimized for the specific device (machine-readable) and is converted from the operating system's format (human-readable). In addition, tapes and disks are usually accessed in physical blocks, rather than a byte at a time. Block sizes may range from 512 bytes to 4K or larger. Each of the devices takes a human-readable file and breaks it into a small piece of the data (determined by the block size and the device controller for either the tape drive or disk) and writes it as a "blob" of data. This data blob is typically sized to the controller's capability, and formatted in a fashion that is convenient for the data interface, say SATA, SCSI or PCIe. This blob of data also contains additional ancillary data for error correction and is encoded in a fashion that is optimal for the device. In short, it will NOT be recognizable as either the PDF document that the user just viewed or the image captured by the camera. To bridge the physical data interface to the computer and to enable the computer's operating system (or kernel) to present the data to users in the notion of file names it is necessary to have a file system. The file system is used by an operating system to organize and manage files on a storage device and defines how data is stored, accessed, and organized on the storage device.

File systems can be viewed as a layered design, as shown in Figure 1:

At the lowest layer are the physical **devices**, consisting of the
magnetic media, motors and controls, and the electronics
connected to them and controlling them.

**I/O Control** consists of **device drivers**, which communicate
with the devices through the controllers at an assembly
language level.

The **basic file system** level works directly with the device
drivers in terms of retrieving and storing raw blocks of data,
without any consideration for what is in each block.
Depending on the system, blocks may be referred to with a
single block number, (e.g., block # 234234), or with head-
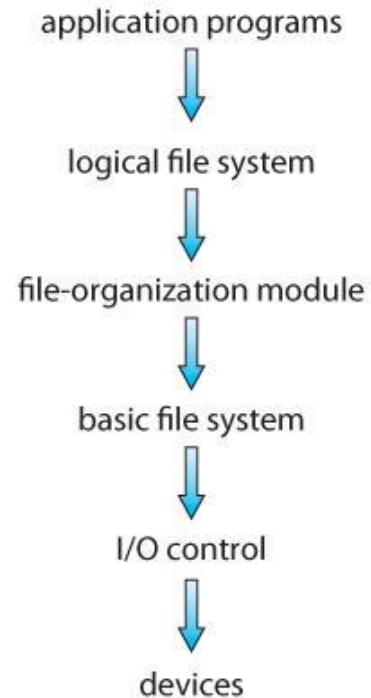sector-cylinder combinations.

The **file organization module** knows about files and their
logical blocks, and how they map to physical blocks on the
disk. In addition to translating from logical to physical
blocks, via a bitmap, the file organization module also
maintains the list of free blocks, and allocates free blocks to
files as needed.

The **logical file system** deals with all of the metadata
associated with a file (UID, GID, mode, dates, etc.), i.e,
everything about the file except the data itself. This level manages the directory structure and the
mapping of file names to **file control blocks, FCBs**, which contain all of the metadata as well as
block number information for finding the data on the disk.

And finally, the **application programs** permit the user to navigate directories and folders, and
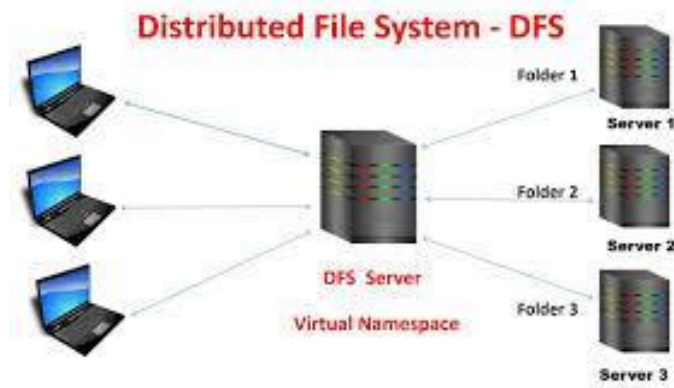save and open files.

On personal computers, file systems began by addressing a single floppy disk and then grew to
addressing hard disks using the FAT16, FAT 32 and NTFS on Windows OS, HFS and HFS+ on
MacOS and ext3 and ext4 on Linux platforms. Similarly, large mainframe computers offered
"access methods" to access data on disk, tape or other external devices. Access methods
provided an application programming interface (API) for programmers to transfer data to or from
a device, and could be compared to device drivers in non-mainframe operating systems, but
typically provide a greater level of functionality.

As storage needs increased, it became necessary for file systems to accommodate greater and
greater pools of storage. Even as single device storage sizes grew beyond the initial 3.75 MB
hard drive, users needed to access gigabytes and terabytes of storage, and this now meant that
storage devices had to be pooled together, and presented to the user in a homogenous fashion –
as though it as one large drive (C: drive) they were accessing. Additionally, with networked

computers, it became necessary to support multiple users, in multiple locations, across multiple operating systems and versions of software. This gave rise to the distributed file system.

A distributed file system (DFS) is a file system that allows one to group shared folders located on different servers into one or more logically structured namespaces. The main purpose of the DFS is to allow users of physically distributed systems to share their data and resources by using a common file system. A collection of workstations and mainframes connected by a Local Area Network (LAN) is a configuration on DFS. DFS also enables spanning across multiple file servers or multiple locations, such as file servers that are situated in different physical places. In a DFS, files are accessible just as if they were stored locally, from any device and from anywhere on the network. A DFS makes it convenient to share information and files among users on a network in a controlled and authorized way.

DFS features, now de facto requirements of any file system, are:

### 3.1. Transparency

The client does not need to know about the number or locations of file servers and the storage devices, and local and remote files shall be accessible in the same manner. Both local and remote files shall be accessible in the same manner, with no hint in the name of the file to the location of the file. Once a name is given to the file, it shall not be changed during transferring from one node to another. Similarly, file naming, when superseded by identifiers and resolution services (see section 4.3 Object Storage), shall not be changed when moved.

### 3.2. User Mobility

A user's home directory shall be directly available on the node where the user logs in.

### 3.3. Performance, High Availability and Scalability

Performance shall be similar to that of a centralized file system. Performance is based on the average amount of time needed to resolve the client's requests. This time covers the CPU time, the time taken to access secondary storage and the network access time.

A DFS shall be able to continue in case of any partial failures like a link failure, a node failure, or a storage drive crash. A highly available and adaptable distributed file system shall have different, independent and redundant file servers for controlling different and independent storage devices.

Since growing the network by adding new machines or joining two networks together is routine, the distributed system will inevitably grow over time. As a result, a DFS shall be built to scale quickly as the number of nodes and users in the system grows. Service should not be substantially disrupted as the number of nodes and users grows.

### 3.4. Data Integrity and Security

Multiple users frequently share a file system. The integrity of data saved in a shared file shall be guaranteed by the file system. That is, concurrent access requests from many users who are competing for access to the same file shall be correctly synchronized using a concurrency control method.

A distributed file system shall be secure so that its users may trust that their data will be kept private. To safeguard the information contained in the file system from unwanted and unauthorized access, security mechanisms shall be implemented.

## 4. Realizing data on storage – File, Block and Object Storage

### 4.1. File Storage

Common implementations of distributed file systems are the Network File System (NFS) and Server Message Block (SMB) and its variants. NFS, a client-server architecture, allows a computer user to view, store, and update files remotely. The protocol of NFS is one of the several distributed file system standards for Network-Attached Storage (NAS). Invented by IBM, SMB is a protocol that allows computers to perform read and write operations on files to a remote host over a Local Area Network (LAN). The Common Internet File System (CIFS) is a variant of SMB and is primarily deployed on the Windows operating system. In this type of storage, known as **File storage**, a hierarchical structure is employed, where files are organized by the user in folders and subfolders, making it easier to find and manage files. To access a file, the user selects or enters the path for the file, which includes the sub-directories and file name.

While this type of storage is easy to navigate, and can store a vast array of complex and different files, the virtual filing cabinet has many drawbacks. One limitation is that accessing and managing data within a file storage array becomes a major inconvenience as the volume grows. The more files, folders, and directories there are, the more time one spends finding and accessing a piece of the required information. Second, the only way to scale out file storage is by adding new storage systems, which can get expensive as a business grows. Third, concurrent usage can pose problems, and while file locking schemes are deployed, they permit two or more users to view a file simultaneously, but are unable to handle simultaneous updates. And finally, many other limitations related to data redundancy, protection and response time.

## 4.2. Block Storage

A solution to the limitations of file storage is **Block storage**. Block storage is where the data is split into fixed blocks of data and then stored separately with unique identifiers. With a unique ID, the storage system is now able to place the data wherever it is most convenient. The blocks can be stored in different environments, such as one block in Windows and the rest in Linux.

This decoupling of data from the user's environment and distribution across multiple environments provides faster and more reliable write access to the data. For reading, when data is requested, the underlying storage software reassembles the blocks of data from these environments and presents them back to the user. It is usually deployed in storage-area network (SAN) environments and must be tied to a functioning server. Since block storage doesn't rely on a single path to data, like file storage does, it can be retrieved quickly. Each block lives on its own and can be partitioned so it can be accessed in a different operating system, which gives the user complete freedom to configure their data. Access to block storage is granted via high performance protocols like Fiber Channel over Ethernet (FCoE) or Internet Small Computer Systems Interface (iSCSI). It is ideal for high-performance mission-critical applications and can provide high I/O performance and low latency.

It's an efficient and reliable way to store data and is easy to use and manage. It works well with enterprises performing large transactions and is the default storage for HDDs and SSDs, however, it has its drawbacks too. Block storage can be expensive, has limited capability to handle metadata, and server binding issues that limit the number of servers that can access it simultaneously.

## 4.3. Object Storage

An efficient and alternative option to block storage is **Object storage**. Object storage is very inexpensive compared to block storage and because it is easy to scale, organizations can use it to store massive volumes of data inexpensively.

On an object store, data is divided into separate, self-contained units that are stored in a structurally flat data environment, with all objects at the same level. There are no folders or sub-directories or complex hierarchies as in file storage. Each object is a simple, self-contained repository that includes the data, metadata and a unique identifying number (instead of a file name and file path). Users can set the value for fixed-key metadata with object storage, or they can create both the key and value for custom metadata associated with an object. This information enables an application to locate and access the object.

Object storage devices can be aggregated into larger storage pools and distributed across locations. This allows for unlimited scale, as well as improved data resiliency and disaster recovery. Objects can be stored locally on hard drives, but most often reside on cloud servers, with accessibility from anywhere in the world. However, unlike with file storage, an Application Programming Interface (API) is used to access and manage objects. The native API for object storage is an HTTP-based RESTful API (also known as a RESTful web

service). These APIs query an object's metadata to locate the wanted object (data) via the internet from anywhere, on any device.

Unlike block storage, object storage doesn't permit editing one part of a file. Objects are considered complete units and can only be viewed, updated, and rewritten as entire objects. Depending on object size, this can negatively affect performance. Another important difference is that operating systems can access block storage directly as attached disks, but they cannot directly access object storage (or if they do, performance suffers significantly). On the other hand, object storage requires virtually no management, unlike block storage which requires remapping volumes and other ongoing maintenance by administrators.

## 4.4. Error Reduction Methodologies

As outlined in *Digital Storage Considerations and Devices for Archiving Motion Picture Assets*, the actual failure rate of the underlying devices of storage, i.e., tape, HDD and SDD is 0.00625%, 1.7%, 1.05% respectively. When combined into large enterprise systems, additional strategies such as checksumming, scrubbing, RAID levels, erasure coding, and data replication need to be deployed to increase the reliability of the storage.

To mitigate data corruption during transfer or at rest, a checksum is computed on the data, and this checksum is recomputed at every stage of transit or when the data is at rest as a verification process. The checksum is a small-sized block of data derived from the original block of digital data for the purpose of detecting errors. Checksums are generated using a checksum function, which outputs a significantly different value, even for small changes made to the input. If the computed checksum for the current data input matches the stored value of a previously computed checksum, there is a very high probability the data has not been accidentally altered or corrupted. Common checksum algorithms include:
- MD5 (Message Digest Algorithm 5): Produces a 128-bit hash value.
- SHA-1 (Secure Hash Algorithm 1): Produces a 160-bit hash value. SHA-256, SHA-384, and SHA-512: Part of the SHA-2 family, these algorithms produce hash values of 256, 384, and 512 bits respectively. They are widely used and more secure than MD5 and SHA-1.
- CRC (Cyclic Redundancy Check): A family of algorithms that produce a checksum, often used in network communications and storage systems.

Additionally, for digital archives, periodic integrity checks are performed at the storage media level, by reading the stored data, recomputing the checksum and validating it against the previous checksum that was computed when the data was written to the media.

Data scrubbing is a process of inspecting volumes and modifying the detected inconsistencies. As time goes by, some data may fall victim to slow degradation that gradually deteriorates data integrity. File system data scrubbing employs the checksum mechanism to check the volumes in the file system. If any data that is inconsistent with the checksum is detected, the system will try to use the redundant copy to repair the data. The file system calculates a checksum for every written file, and further protects that data checksum with another checksum (metadata checksum). Every time data scrubbing is

conducted, the file system recalculates the checksum and compares it with the previously stored data checksum. Additionally, the data checksum will cross-check its corresponding metadata checksum to make sure the data checksum itself is intact. Once data corruption is detected, the system will try to repair the corrupt data by retrieving the redundant copy
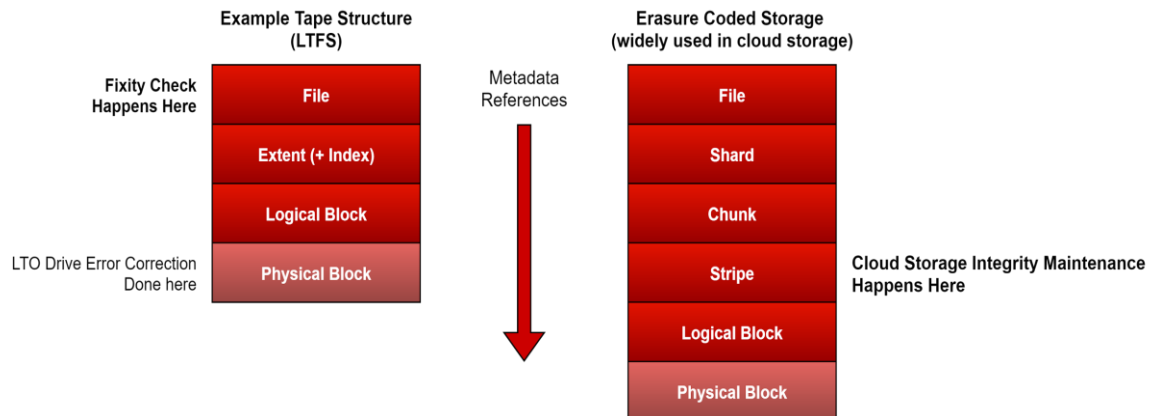
RAID stands for Redundant Array of Independent Disks. Simply put, it combines multiple drives into a single storage pool, thus offering fault tolerance and data redundancy. Many RAID levels employ an error protection scheme called "parity," a widely used method in information technology to provide fault tolerance in a dataset. RAID 5 is a commonly used methodology and requires at least three drives and utilizes parity striping at the block level. In distributed systems with three-way replication for data protection, the original data is written in full to three different drives and any one drive is able to repair or read the original data. While this replication provides data recoverability, it is not the most efficient use of expensive storage, and suffers from the fact that on drive failure the system is placed in read-only mode at reduced performance while data is copied onto a new drive to replace the failed drive.

In object storage, a common error mitigation scheme is erasure coding. Erasure coding splits data files into data and parity blocks and encodes it so that the primary data is recoverable even if part of the encoded data is not available. Object stores utilize erasure coding to provide data protection by saving encoded data across multiple drives and nodes. If a drive or node fails or data becomes corrupted, the original data is reconstructed from the objects saved on other drives and nodes. Erasure coding is able to tolerate the same number of drive failures as other technologies with much better efficiency by striping data across nodes and drives. For example, 10 PB of data replication would require more than 30 PB of storage in RAID 5, whereas object storage would only require 15-20 PB to securely store and protect the same data using erasure coding.

### 4.4.1. A Note on Fixity

Fixity[2], in the preservation sense, means the assurance that a digital file has remained unchanged, i.e., fixed. In digital archives, checksums are extensively used to verify the fixity (integrity) of files stored on disk or when transferred across devices or storage media. By calculating a checksum for a file and comparing it to a known good checksum, users verify that the file has not been altered or corrupted.

While fixity is done at the file system level, on object stores data integrity checking is done at a lower abstraction layer in the storage stack, a level at which the data integrity algorithm has no visibility into filesystem structure. See Figure 3 below – (from the ETC paper on cloud archive[3]). Hence, on object store, the notion of fixity may not have as much relevance, and the only mechanism to carry out fixity checks on object store is to access the archive files (which will be visible as objects, since this is the layer at which they operate) and perform the hash calculation and comparison.

**Example Tape Structure (LTFS)**

Fixity Check Happens Here — File
Extent (+ Index)
Logical Block
LTO Drive Error Correction Done here — Physical Block

Metadata References

**Erasure Coded Storage (widely used in cloud storage)**

File
Shard
Chunk
Stripe — Cloud Storage Integrity Maintenance Happens Here
Logical Block
Physical Block

Fixity checks and reports, as pointed out in the ETC's Practical Cloud Archive[3] paper, have "become an essential fiduciary requirement of archive management, in particular for owners of extensive collections of intellectual property." However, the paper also recognizes that "the current fixity verification implementations are based on the current wisdom and experience of the durability of the underlying storage medium - LTO tape" and that newer storage like object storage will have to be studied to better understand its durability and if the mechanism of fixity as applied to tape is valid.

## 5. Cloud

Cloud storage is a remote storage model that enables storing data at a physically remote site and accessing the remote computing provider either through the public internet or a dedicated private network connection. The provider securely stores, manages, and maintains the storage servers, infrastructure and network to ensure access to the data when you need it at virtually unlimited scale, and with elastic capacity. Cloud storage removes the need to buy and manage one's own data storage infrastructure.

Such a storage model alleviates the need to buy and manage one's own data storage infrastructure, providing agility, scalability, and durability, with anytime, anywhere data access.

Cloud storage is delivered by a cloud services provider, and this could be a public provider like AWS, Microsoft Azure, Google, etc., or a private cloud provider like Wasabi, Seagate, etc. The cloud provider owns and operates large data centers in multiple locations around the world, and manages capacity, security, and durability to make data accessible to the user's applications over the internet in a subscription based (pay-as-you-go) model. Typically, the storage cloud is connected to the user either through the internet or through a dedicated private connection, using a web portal, website, or a mobile app. Applications access cloud storage through traditional storage protocols or directly using an application programming interface (API).

Cloud providers typically offer different tiers of storage for their customers which are based on access patterns (e.g., daily, monthly, quarterly, yearly, long-term). Each of these tiers has a different pricing model, and in many cases, a different API is used for access. Additionally, each tier typically has an early retrieval fee if the content is accessed before the retention period of that tier has passed. In some cases the least expensive archive tier uses LTO tape for storage of the files (e.g., AWS, Microsoft, IBM), and in others the coldest archive tier is spinning disc (e.g., Google).

All three types of storage discussed in Section 4 are offered by cloud providers, and depending on the user's performance needs and budget, compute power and data capacity are easily scaled for file, block or object storage. As an example, the commercial offering from Amazon, AWS provides:

- Object storage as the Simple Storage Service (S3) where objects are stored on multiple devices across a minimum of three discrete data centers (Availability Zones).

- File storage as Amazon FSx and Elastic File System (EFS) where shared file access for applications with unstructured data such as audio, video, images, and user directories are provided.

- Block-based storage as Elastic Block Store (EBS) to deliver ultra-low latency for high-performance workloads.

Cloud services are very secure, and permit fine-grained control on where the data is stored, who can access it, and what resources can consume it at any given moment. Fine-grain identity and access controls combined with continual logging and monitoring ensures information is only accessed by those with the right access permissions.

For cloud storage, durability and availability are the two main metrics that are used to quantify data protection and access guarantees.

- Durability is the metric that quantifies the safety of the data from irrecoverable loss, specified as a percentage. Currently, services such as Amazon S3, Microsoft Azure and Google cloud claim to provide a durability of 11-nines (99.999999999%) or $10^{-11}$ over a given year. This durability number means that even with one billion objects, one would likely go a hundred years without losing a single one.

- Availability metric quantifies the amount of time for which data is available for access. Typical availability figures for cloud storage services are around 99.9%-99.99%, which translates to approximately 10 hours to 1 hour of unavailability in a year respectively.

## 6. Value-Added by the Cloud

Cloud providers provide more than just storage. Cloud infrastructure also offers capabilities such as compute, networking, artificial intelligence/machine learning and security, and these form the

basic building blocks of a holistic cloud solution. Additionally, for archival purposes these services can be leveraged and scaled beyond what could typically be achieved in an on-prem environment, and these services can be expanded or contracted as required by the use case. One specific case, given the accessibility of compute resources and the proximity of the storage is that traditional fixity reports can be generated by performing checksum calculations – an essential component of fixity checks – on the objects that have been archived.

Beyond redundancy and elasticity in infrastructure, specific solutions that are designed to address specific needs can be deployed as Software-as-a-Service (SaaS) on the cloud and these applications act as the interface between the user and the underlying cloud resources. An example of a media industry-specific solution is a Media Asset Management system. A cloud-based media archive would typically include a MAM system involving a number of different infrastructure solutions (e.g., compute, storage, networking and security) in combination with application-specific solutions such as the MAM software, transcoder applications, databases, and specific software applications.

Digital archives can benefit from these specific solutions, as they can fully address the enterprise-level needs of ontology, metadata management, policy enforcement, governance, etc.


Listed below are a few (and by no means all) of the solutions that a cloud-based archive could benefit from, given the elastic compute capabilities, tiered storage, geographical replication, and software services that can be deployed.


- Image and sound metadata enrichment

Currently the majority of the metadata tags are curated by humans. Utilizing the large pool of compute resources and large language models, artificial intelligence and machine learning (AI/ML) tools can generate descriptive metadata tags for media, describing both the image and sound details of the content. Additionally, relevant metadata can be scraped from already available sources such as camera reports, script supervisor notes, location information, sound mix logs, call sheets, etc. Future-proofing metadata requires a lexicon that can be difficult if not impossible to maintain and update.

In addition to descriptive metadata, the metadata tags serve as a more practical approach to use for search. With enriched metadata tags, it is then possible to search and identify content in a manner that doesn't require knowledge of the actual data or search terms previously defined.

- Image and sound data ontology

The definition, tagging and linkage of media files helps organization and search of the archive. However, the success of the search is predicated on not only the accuracy of the data, but also on the use of very specific search terms that describe the contents as well as the linkage between them. These search terms may not be obvious to future generations of archivists, which can cause difficulties in finding

specific pieces of content. It is therefore important to build data ontologies as a method for organizing and structuring data based on the relationships between content files. In a storage structure where individual objects are referenced (as in an object store), numerous relationships can be established, either from a versioning standpoint or from a generational standpoint, or even from a genealogical standpoint. Similarly, ancillary data such as camera reports, script supervisor notes, location information, sound mix logs, call sheets, etc. can be linked directly to the media and this relationship can be maintained throughout the life of the archival elements.

- MAM systems

Media asset management (MAM) systems enable media companies to manage media on the cloud. Critical components of such systems are

- Centralized uploading of content with permission-level access
- Tagging of video content based on enterprise-level rulesets with learning capabilities
- Proxy and original content downloads as needed
- Collaborative editing for video and graphics productions
- Centralized storage location for easy access

With centralized storage, content owners can consolidate their data and maintain ownership of it, rather than have this managed by service providers. This gives them complete control over their content. With cloud services, MAMs enable producers, editors, and artists to collaboratively review, edit and approve creative decisions on their productions. This greatly simplifies the production process and enables a faster time to market.

- Digital Distribution

  Digital archives can (and should) become a central part of a media company's digital supply chain, supporting the delivery of content to clients and customers. With the elastic compute capabilities to transcode the audio and video assets to any downstream distribution specification, the cloud-resident archive can support many direct-to-consumer streaming platforms in addition to traditional customers such as VOD platforms, theaters, cable companies, broadcasters, disc-based retail, etc.

- Image and sound restoration applications

  Archivists have always dealt with archiving the content as it is presented. In many cases it may be damaged and need restoration.

With a consolidated archive in the cloud, and the correct ontology it is possible to identify elements that may have higher quality and greater fidelity.

It is also easy to deploy digital restoration techniques that have the benefit of AI/ML and apply additional heuristics to enrich the digital archive.

- Language Dubbing and Subtitling

A motion picture's worldwide release footprint and schedule dictate the number of different languages that the audio is dubbed into, as well as the different languages and dialects of subtitles that are generated. Without a consolidated and collaborative platform like the cloud, spoiled low-resolution versions of an early cut are digitally shared (sneaker net, or over the internet) with individual dubbing studios and subtitling facilities. Once dubbing is completed, the audio files and subtitles are returned to the production, which has progressed forward with additional editorial changes, which now in turn have to be re-dubbed. In a centralized store, dubbing studios can be on-boarded much faster, and all changes can be monitored by all parties, bringing numerous efficiencies to the localization effort.

Further, with the advent of AI/ML tools, content owners are adopting cloud services to carry out automated captioning[4] and many are investigating cloud services for automated language dubbing.

- Image and sound quality control applications (QC/QA)

Media and entertainment assets are among the most scrutinized from a quality assurance and quality control (QA/QC) standpoint. Every frame of a master file is visually inspected, and done at every hand-off stage. Similarly, after every transfer to a post-production facility or whenever a copy is made, QC is carried out at the same level of diligence. By eliminating the need to transfer or make copies, and directly interact with a single source, the need for repeated QC is eliminated. Further, by establishing the right ontology, the master data's QC report can be correlated with any derivative version. This way issues that are identified (and rectified) in the master need not be flagged in downstream versions. Automated QC tools can generate test reports that are hierarchical so that a historical sequence of identified errors and fixes can be recorded.

Further, with the advent of AI/ML tools, content owners are adopting cloud services to carry out automated QC and many are investigating the use of the result of the automated QC to determine downstream processes such as transcoding and packaging.

- Deduplication of assets

Bringing all the different versions, language versions, aspect ratios and other variants of a motion picture together in one place enables digital

archivists to examine all the places where the content is duplicated, and if necessary move the duplicated content off to much cheaper storage or, depending on the level of comfort in the deduplication algorithm, delete the content, and recover storage.

This can be applied to localized content and distribution content very easily, as well as extended to the multiple takes that are captured during original photography.

- Disaster Recovery and Business Continuity

All cloud providers provide different replication strategies to ensure data recovery in the event of a disaster or even business continuity in the event of an outage. In addition to replicating the content within a specific region, to protect against such catastrophic natural events, it is recommended to replicate the data across geographical locations. In such scenarios, the replicas of the data are stored in different geographical locations of the world, enabling data access in the case of a hurricane or earthquake in a specific region.

Additionally, with the compute and network infrastructure offered by the cloud providers, with the content stored in the cloud, a media company can continue to serve its digital supply chain in the event of a catastrophic event.

- Security

Cloud storage not only provides robust security features like file encryption at rest and in transit, and secure access controls to safeguard assets, it also ensures that the data storage complies with various laws and regulations. This means that in addition to the data being safe from unauthorized access it is also stored in a way that meets legal standards.

As described above, access to cloud data and services are through APIs, and all API services should require individual Identity Access Management and authentication. Additionally, with the use of virtual servers, web application firewalls, and robust Cloud Security Posture Management[5], all the leading cloud providers have aligned themselves with most of the well-known accreditation programs such as PCI 3.2, NIST 800-53, HIPAA, GDPR and the TPN security requirements of the Motion Picture Association.

- And many other use cases that are being deployed as we speak, such as Marketing, as are concepts such as the elements in a "digital backlot"/digital twins and other asset reuse cases for games, immersive applications, etc.
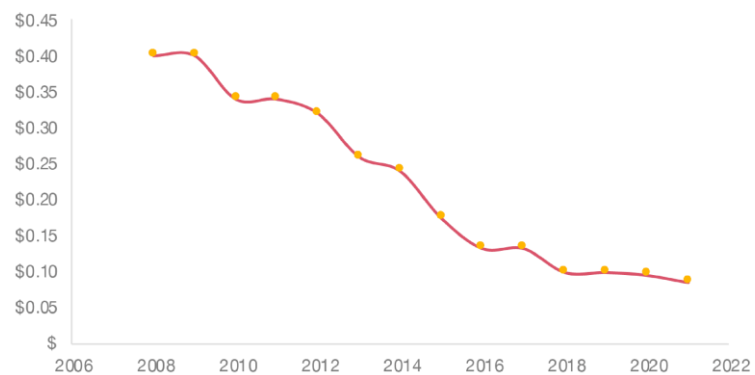
The Academy Digital Preservation Forum encourages readers to submit their utilization of the cloud for their digital archives in the form of a post in the forum.

# 7. M&E Challenges for Cloud Archiving

As seen above, there are numerous advantages to adopting the cloud; either private or public. While the term "cloud" may seem foreign, Sections 3 and 4 illustrate that the cloud consists of standard storage devices and relies on the same devices that digital storage is stored on today. Given these benefits, companies are increasingly adopting the cloud, and per Gartner the spend on cloud services is forecast to grow 20.4% to total $675.4 billion in 2024, up from $561 billion in 2023 with

- 96% of companies using at least one public cloud.
- Companies running 50% of their workloads in a public cloud.
- Organizations storing 48% of data in a public cloud.

At the same time, the cost of cloud computing is tracking to "Bezos' Law" (named after Jeff Bezos, founder of Amazon), which states that the cost of cloud computing will be cut in half every 18 months – as can be seen in Figure 4, which shows the USD cost per hour for Linux on-demand from the us-east-1 region on AWS.



There are, however, challenges to cloud adoption, across all industries and some specific to media and entertainment

- Costs

Though cloud migration brings in a lot of returns and benefits, in the long run, getting there is usually expensive and time-consuming. Costs include architecture changes, human resources, training, migration partners, cloud provider, and bandwidth costs.

Additionally, these have to be balanced with the cost of maintaining a local infrastructure, as most media companies already have some legacy on-prem infrastructure.

- Asset Curation and Migration choice

Deciding on what to migrate to the cloud or to archive in the cloud is vitally important for any archive. Therefore, careful selection is extremely important. Additionally, it is important that the assets are carefully curated to ensure that the dataset within the archive is clean. Many times, as the assets may be at different service providers, the asset

curation process can be time consuming and frustrating.

- Complexity and Lack of Expertise

Most organizations are wary of and inexperienced with the complexity of cloud environments and migration processes. Most media companies rely on service providers to fulfill parts of the production chain, as well as fulfill the deliveries in the digital supply chain. With cloud environments the expectation is that the curation and fulfillment can be brought in-house, and if there is not enough in-house expertise in dealing with cloud, migration processes, and compliance requirements then adoption is likely to be slow.

- Cloud vendor lock-in

The availability of multiple similar cloud service providers makes it a hurdle to choose the right one. Goals, budget, priorities of the organization, along with the services offered, security, compliance, manageability, cost, etc. of the service provider should be the main considerations when making a selection.

However, once a provider is selected, the utilization of services on that provider could force a content owner to feel locked-in to that provider, as not all services are provided on all cloud platforms, at the same quality of service. Further, the modifications needed to operate on another provider may prove to be inordinately expensive.

Further, public cloud vendors still operate on a model where there are additional charges for egress, so making a cloud selection can have very significant cost implications, especially if at some point data has to be transferred out. There are cloud providers who currently do not charge egress fees, and it will be interesting to see if the public cloud vendors offer to waive egress fees also.

Another concern with a single vendor is that once all data has been migrated to the cloud, the media company will be bound to the same vendor even if their prices do not remain competitive with respect to other providers. This is especially important given the steadily decreasing costs of cloud services globally. With these constantly falling prices, the last thing a company wants is to be tied to an uncompetitive vendor.

- Interoperability between cloud systems

    Not all services are equally provided on all cloud providers. While some software applications are supported on multiple clouds, they typically tend to be optimized for a particular cloud vendor. This is for many reasons, the main one being the level of engineering needed. While at a base level, each of the cloud providers provide similar functionality, to take full advantage of any one cloud provider it is necessary to integrate many custom features. While these features may

be available on another cloud provider, the level of effort and engineering complexity can be fairly high to fully take advantage of all of the optimizations available.

Cloud APIs that transfer data between cloud computing services or between cloud services and on-premise applications are currently available, but they may not be compatible with every cloud provider or even be designed to work across different providers' environments.

- Non-Digital Assets

  Media companies have been storing physical assets for a long time, either on local vaults or offsite. Any holistic archiving system has to be able to consolidate digital and non-digital assets and present a complete view of the archive to the user, so that meaningful decisions can be made regarding reuse, restoration and monetization. There are very few commercial systems available that can support this dual functionality out-of-the-box.

- Data Security

  Data security is the biggest concern when enterprises store their sensitive data with a third-party cloud provider. If data is lost, leaked, or exposed, it could cause severe disruption and damage to the business.

## 8. Conclusion

Private or public cloud storage is built upon storage that is commercially available and has been deployed on-prem and in data centers for quite a while. The addition of cloud computing makes it attractive for archiving, as it can offer several benefits:

- Potential cost savings from easier procurement and economies of scale, particularly for smaller repositories, and if using the lowest cost tier of storage.

- Cloud services can provide easy, automated replication to multiple locations, essential for disaster recovery and business continuity for a content owner.

- Specialized software can be easily deployed providing access to dedicated tools, procedures, workflow and service agreements, tailored for digital preservation requirements.

- Easy to pilot emerging service providers and tools and carry out rapid and low-cost testing within a Software-as-a-Service (SaaS) model.

With the increasing commercial availability of AI/ML tools, the combination of large compute pools provided by a cloud provider and the large datasets presented by the consolidated archive of a media company will almost certainly result in numerous additional use cases for the transfer of data to the cloud, either private or public.

Additionally, the cost of cloud services and infrastructure continues to drop, and is tracking to a predicted rate of halving every 18 months, so what was expensive yesterday may no longer be today.

Given the continued expansion of the digital storage footprint of motion pictures and the increasing use of media creation cloud workflows, the need for secure, collaborative workflows, and shorter delivery times, it is abundantly clear that there exists a need for secure, reliable, consolidated storage with on-demand compute and networking infrastructure – a need that is fulfilled by cloud providers. However, as illustrated in Section 7, there are unique challenges posed by the motion picture industry on the cloud; however, these are not insurmountable, especially when approached with the correct strategy and fallback options. One option, pointed out in the ETC[3] paper, is to not view cloud storage in an archival solution as an all-or-nothing scenario. Another approach would be to use the cloud as one of the archive elements in a multi-element approach and instead of making three tape copies, have two tape copies and one cloud copy. This will permit a content owner to examine the value-added services illustrated in Section 6, and gain confidence in their cloud provider.

## 9. References

1. https://www.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/12_FileSystemImplementation.html
2. http://blogs.loc.gov/thesignal/2014/04/protect-your-data-file-fixity-and-data-integrity/
3. https://www.etcenter.org/wp-content/uploads/2022/02/Practical-Cloud-Archive-FINAL-02082022.pdf
4. https://www.hollywoodreporter.com/business/business-news/warner-bros-discovery-google-captioning-1236010573/
5. https://www.microsoft.com/en-us/security/business/security-101/what-is-cspm
6. https://www.gartner.com/en/newsroom/press-releases/2024-05-20-gartner-forecasts-worldwide-public-cloud-end-user-spending-to-surpass-675-billion-in-2024