

ENTERPRISE FLIGHT DATA MANAGEMENT SYSTEM (EFDMS) AND STORAGE INFRASTRUCTURE TECHNOLOGY DISCUSSION

**Robert Crenwelge, Consulting Systems Engineer
EMC Corporation
Austin, Texas**

**Brian Conway, Professional Services
EMC Corporation
Edwards Air Force Base, California**

**Kevin Dillon, Centera Senior Systems Engineer
EMC Corporation
McLean, Virginia**

ABSTRACT

This Paper presents efforts in developing a data management system and storage infrastructure for assisting test engineers in achieving information superiority and maintaining vital up-to-date information. The focus of this Paper is to generate support for a technology refresh, upgrading the major data centers that share in the responsibility of processing telemetry information. We illustrate how our efforts fit into this goal and provide an overview of our concept for a revolutionary transformation in data management systems. We present the significance of this new technology and suggest a path to implementing the solution.

KEY WORDS

Flight Test, Telemetry, Data Storage, Archival, and Networks.

INTRODUCTION

Since the 1980s, test engineers have faced growing challenges in information management and data processing. The number of measurement parameters and sampling rates has skyrocketed with every new generation of aircraft and test article. This creates enormous amounts of data, and as aircraft become more sophisticated, the amounts of data will increase exponentially. Since its inception, telemetry data for our Nation's premier fighter aircraft has been captured on tape and subjected to a labor-intensive process for analysis. The telemetry data is then requested by authorized users and satisfied by shipping duplicate tapes to the requestor.

The Enterprise Flight Data Management System (EFDMS) addresses the issues that surround storing, managing, retrieving, and distributing telemetry data. Data can be downloaded from a solid-state memory canister, loaded to the EFDMS, and accessed instantaneously. Security is a top priority for protecting sensitive information and with the robust features of commercially available applications, sensitive or classified data can be protected with policy-based features granting access to data by organization, group, or at an individual level. By simply automating manual processes, the collection and dissemination of data to the engineer can be streamlined resulting in 24x7 access to mission-critical data.

The increase in the quantity and quality of the various types data will contribute to a higher quality finished product that can be delivered in record time. New types of data will begin to appear as capacities increase and engineers begin to mix telemetry, voice, and video to have more information available at their desktops.

The EFDMS concept can provide years of extended capabilities to fixed content, whether the data consists of telemetry, video, voice, or a future data source. Instantaneous access to information will be a determining factor in maintaining information superiority.

The intent of this paper is to generate support for a technology refresh upgrading the major data centers that share in the responsibility of processing telemetry information. The benefits achieved by eliminating inefficiencies needs to be explored in detail with a consolidated cost recovery analysis.

UNDERSTANDING THE PROBLEM

Performing testing maneuvers of fighter aircraft can wreak havoc on conventional tape drive technology. It becomes increasingly more difficult to collect data for analysis when multiple flights are required to collect sufficient data.

Solid-state devices are currently being tested as a replacement for conventional tape, promising higher levels of availability, reliability, and scalability and have proven to be successful (i.e., resulting in no loss of data). These devices provided a tremendous increase in quality and quantity of data produced for flight test engineers. However, today's data and information processing requirements have outpaced the capacities and capabilities of these legacy management procedures.

The current procedures and policies for data archival and day-to-day management of telemetry data is labor intensive in every aspect. As the manpower resources available to perform these tasks has decreased, new solutions must be considered to provide the capabilities to manage, protect, and share mission-critical information. Additionally, with the programs such as the Joint Strike Fighter (JSF) in place, and the continued need to leverage new weapons technology with older aircraft, the corresponding demand to record, access, and process telemetry data is even more crucial.

From the time a request is received, processed, and then shipped, days can pass before the end user actually receives the data. Each data request is satisfied by creating duplicate tapes, a process which potentially reduces the useful life of the master tape. This could place the master tape at greater risk for data loss. The subsequent shipment of the tapes to the end user is often subject to security

procedures that could jeopardize the timely receipt of user data. This workflow process can often take hours and sometimes days before the engineers are able to load and process the data for analysis. Furthermore, this process can create a significant backlog of data requests that can potentially cause project delays, security breaches, and unavoidable costs.

As noted previously, the current method for processing telemetry data from aircraft to end-user is cumbersome. The following scenario describes a typical, time-consuming, step-by-step process from data generation to collection to end user analysis –

1. Once data has been recorded, the solid-state memory canister is removed from the aircraft and taken to a workstation, to which the solid-state memory canister is then attached.
2. The data is then downloaded to a workstation. At this time, the data is recorded to tape for archival and sometimes analyzed by a local engineer.
3. The original tapes are then placed into a vault for long-term storage.
4. The data does not live anywhere on disk or in a Storage Area Network (SAN) for instant access and retrieval.
5. At some point, an engineer determines a need to analyze test flight data from the flight or from a series of flights several months back.
6. The engineer must send a request for the data to the staff managing the tape archives.
7. The staff must then research a lookup table to find the tape(s), which contain the requested information.
8. The tapes are then pulled for processing. A resource is then assigned to duplicate the tapes, which can require hours to complete and depends heavily on the amount of data requested.
9. Before the tapes can be sent, a security clearance check must be completed to verify that the recipient, shipping location, etc., are approved for taking delivery of the sensitive or classified data.
10. The tapes are then packaged for delivery through commercial shipping channels, US Postal Service, Federal Express, UPS, etc., or via special courier.
11. It can be expected that several days later, the tapes are received by the requesting organization, and then finally delivered to the engineer who placed the original request.
12. The requesting engineer may then load the data directly onto his workstation, or the engineer may assign the task of loading the data into the organization's storage network to a technician.
13. Once the data has been loaded, the engineer can begin analyzing the data.

These steps demonstrate how difficult and cumbersome the task of managing fixed content for storage, retrieval and distribution can be.

With this in mind, we identified the following operational issues, technical limitations, and goals through our discussion with U.S Air Force personnel –

- Current data duplication processes places mission data at greater risk for loss;
- Access to data is time consuming, inefficient, and costly;
- Current data processing methods place mission data at risk for security breaches;
- The need for automated processes (hardware, software, networks, etc.) is required to address data growth and limited personnel resources within military installations;

- Existing telemetry data is stored on conventional tape, adversely affecting the ability to manage, share, and protect data;
- Solid-state recorders are becoming the solution of choice. Solid-state recording devices can store more than 400 GB outpacing its tape rival by 850%; and,
- The vast amount of data currently being collected can no longer rely on conventional tape technology.

Understanding the problems that telemetry data collection has introduced into flight data management systems is key to designing a new architecture that addresses the downfalls and limitations of tape-based systems and solid-state disk conversions. A new IT storage infrastructure would provide higher levels of data security and integrity, thus eliminating labor delays in data delivery to key personnel or project partners for analysis. This would result in a reduction of engineering cycles and the lowering of overall project costs.

PROPOSED TECHNOLOGY APPROACH

EFDMS provides a flexible infrastructure that allows frictionless access to data. Implementing an infrastructure approach to data collection and dissemination, efficiencies can be realized in performance, information assurance, and addressing scalability requirements. Building an infrastructure for the Flight Data Management System will enable the system to be Enterprise class architecture, one that addresses all aspects of operations including mission continuance, disaster restart, and disaster recovery, all of which are critical components of an enterprise-class data center.

The future IT infrastructure will be the heart of EFDMS's successful operations, providing continuous availability of networks (e.g., SAN or NAS), applications, servers and information storage. In the area of information storage, our approach incorporates a Content Addressed Storage (CAS) solution. The design is targeted at meeting the unique requirements of fixed-content management. Fixed content comprises any form of digitized information asset retained for reference and value, including documents, e-mail, telemetry, satellite photos and imagery, streaming video/audio, X-rays, final form CAD/CAM drawings, etc. Federal civilian agencies and military services that require WORM (write once, read many) storage capabilities will find our approach a viable option. The CAS solution provides fast and easy online access to petabytes of information (PB=1,000TB).

In the new technology technique called content addressing, a unique identifier is assigned to each stored object so specific content can be immediately accessed through an integrated application. To meet the requirements for long-term storage of fixed content, the CAS solution actively ensures that objects cannot be changed or modified once stored. Applications are integrated with the CAS solution via an API (application programming interface) over an IP network, recognizing the CAS as a storage repository, scalable to millions of objects and petabytes of fixed content. Customers and independent software vendors (ISVs) that incorporate the APIs can easily integrate their fixed content applications with the CAS solution.

In our approach to CAS, a flat address scheme is presented to the archive program, as opposed to traditional location-based file methods. When a reference data object is initially stored in the CAS repository, the application is given a "claim check" that is uniquely derived from the object's

content. To access a stored data object, the requesting agent simply provides the “claim check” that uniquely identifies the object to the repository, and the data object is then returned. Content addressing greatly simplifies the storage resource management tasks, especially when handling large amounts of static objects.

Accessing CAS

Our approach synchronously stores a data object, also known as a BLOB (Binary Large Object). An application delivers a data object to our Application Programming Interface (API), which calculates a 128-bit claim check, i.e., Content Address (CA), from the object’s binary representation. The CAS repository then stores the BLOB and a mirror copy. The Content Address, which is unique for the object’s content, and metadata about the object (e.g., filename, creation date, etc.) are then inserted into an XML file called a C-Clip Descriptor File (CDF), which in turn has its content address calculated. In our approach, the C-Clip is the union of the CDF and its content objects. This C-Clip Content Address is only returned to the application once two copies of the CDF and two copies of the BLOB have been safely stored in the repository.

Future requests for the retrieval of the data object occur when an application submits the C-Clip’s Content Address for that object to the repository via the API. There is no centralized directory in CAS and no pathnames or URL’s are used. Only the C-Clip’s Content Address is used as a reference. The C-Clip is essentially a “fingerprint” that assures the authenticity of the stored object (i.e., the user file BLOB). If an object is retrieved and altered by one bit, the CAS API will produce a new CDF with a new Content Address for the altered BLOB when the modified object is subsequently stored. The significance of this approach is that the original object remains unchanged and accessible by its original content address. This WORM (write once, read many) attribute assures a level of versioning integrity that file servers cannot provide.

C-Clip Functionality

The C-Clip method ensures that application developers, users, and storage managers do not need to think about where data is physically or logically located. The C-Clip’s Content Address is a globally unique identifier permitting a data object to be retrieved from anywhere, exactly as it was stored, irrespective of the location of the content and the user.

Because our CAS approach uses a location independent addressing scheme, the result is data mobility that facilitates a simple yet robust disaster recovery topology. When a data object is initially stored in the local CAS, the object can also be asynchronously and automatically replicated to remote sites over a wide area network (WAN), such that the object resides both local and remotely.

If a disaster in the local data center occurs and the data needs to be rebuilt, the remote site is able to re-populate the local CAS as needed. Note that the replication facility can be operated bi-directionally in an active-active topology.

User interface with our implementation of CAS is through content-based software applications that incorporate a powerful API. Many industry-leading software developers have signed up to integrate their products to take advantage of the features that our approach offers. Examples of this type of third-party software that can use our API and CAS include content/document management, medical

imaging, e-mails archiving, and a wide variety of vertical applications that benefit when large amounts of content objects are accessible online and shareable by a large number of users.

Our Architecture

The CAS architecture we implemented presents a *no single point-of-failure* platform that is highly scalable and implements non-disruptive servicing. The CAS solution is built upon a Redundant Array of Independent Nodes (RAIN) that is deployed in one or more six-foot NEMA standard 19-inch racks. A single rack can hold 16, 24, or 32 nodes to provide 4.7TB, 7.2TB or 9.6TB of protected capacity, respectively. Each node contains processing power, 600GB of raw storage capacity, and is interconnected with all the other nodes in the cluster via a private LAN. Each node executes an instance of our software in one of two operational modes:

- Storage Node: the node facilitates long-term storage of BLOBs and CDFs, and
- Access Node: the node is a conduit for interaction between the application server and the storage nodes.

The throughput needs of the application (e.g., telemetry data) determines how many access nodes must be configured at the time of installation. Each access node is connected to the application server infrastructure via a 100 megabit-per-second Ethernet cable. Given that our CAS solution will be connected to the application server via multiple access nodes, the application will enjoy not only scalable bandwidth with low access latency, but also high availability.

This architecture provides our CAS solution extreme scalability for capacity and performance, as well as radically unique simplicity for managing physical storage resources. Capacity scaling is accomplished merely by non-disruptively adding nodes in 2.4TB (protected) chunks. As capacity is scaled, performance is also scaled to handle the management of data across the incremental nodes. Additionally, the number of Access Nodes can also be augmented to facilitate more bandwidth to the applications server.

Perhaps the most significant benefit of our architectural implementation is that the addition of nodes requires no time-consuming, complex management effort on the part of either the application (or developer) or the system administrator. This benefit permits a single administrator to handle hundreds of terabytes of content storage – far greater that can be managed in traditional storage products.

Multiple racks can be configured as a single “cluster” delivering up to 154TB of protected storage. The repository connects directly to one or more Windows or UNIX application servers via multiple TCP/IP LAN connections. If 154TB is insufficient, the client application API can reference several 154TB clusters to access up to one petabyte of data. Building such a network of clusters is simply a matter of non-disruptively connecting new clusters to the application server’s LAN infrastructure as older clusters are filled. The API is designed to abstract the multiple clusters from the application server.

We have begun teaming with key software vendors to provide solutions that eliminate the daunting task of managing fixed content. Today, companies and government agencies are solving the problem of fixed content management with our CAS solution and with software from one of our partners who

have invested time and money integrating their applications to take advantage of the CAS solution's unique strengths and capabilities.

CONCLUSIONS

An EFDMS using our CAS solution for its data repository will completely address all of the existing problems we identified earlier in the process of telemetry data collection, migration, storage/vaulting, and access/retrieval. Manpower, scalability, security, integrity, and disaster recovery issues and problems surrounding these valuable data assets can be virtually eliminated. The CAS solution delivers the following key benefits:

- Simplified management.* Resident applications are not required to understand and manage the physical location of stored information. Instead, the CAS solution creates a unique identifier that applications can use for retrieval.
- Fast, location-independent access.* The CAS solution provides users with fast, shared, and networked access to fixed content at Internet speeds. Access is location-independent, greatly simplifying application development and deployment.
- Assured authenticity, efficient replication.* Because the CAS solution gives each stored object a unique address, integrity and authenticity are assured. Only one copy and one replica of each object is stored, no matter how many times the object is used.
- Complete solutions.* Our CAS solution is already being integrated into many leading content management applications across a variety of industries by an expanding list of partners.
- Scalability without reconfiguration.* The CAS solution architecture is based on redundant arrays of independent nodes (RAIN), making it highly scalable to hold petabytes of content. To add capacity, just plug in another node. Our solution auto-discovers and configures the new capacity as it is installed.
- Self-healing.* The CAS solution continuously monitors to detect and repair soft errors. It also automatically reconfigures itself and replicates objects as necessary, in the event of hardware failures such as disks or nodes. These incidents can be automatically reported through a remote monitoring system.
- Future-proof architecture.* Long-term data often outlives the technology on which it is stored. Because our CAS solution is designed to accommodate new technology without costly and disruptive conversion or migration, it eliminates this problem.
- Business continuity protection.* All information objects are synchronously mirrored within a local CAS cluster to support automatic recovery from component failures. The solution can also be configured to maintain duplicate copies of fixed content at a remote site to guard against local site disasters.
- Easy installation and non-disruptive upgrades.* The CAS systems can be installed or upgraded in less than an hour, without disrupting content access. The CAS solution's software environment can also be upgraded non-disruptively as new versions are released.

Our CAS-based EFDMS solution eliminates tape from the workflow process, removing the most labor-intensive activities from the data center. Although tape and optical solutions may appear less expensive due to lower initial implementation costs, in the long term these solutions require significantly increased focus on manual content management (e.g., movement of tape, conversions to

new formats, copying of old tapes to avoid loss of data due to tape aging effects, etc.). Furthermore, the traditional solutions are more difficult to manage during migrations and changes in technology. The solution presented in this paper also ensures increased productivity when implemented across multiple data centers (e.g., by preventing outages caused by natural disasters), and the data mirroring capabilities in our solution provides for simultaneous data access from multiple data center locations. The savings opportunities associated with the CAS solution has already compelled the commercial sector to invest heavily in this new technology, including EMC's implementation of the CAS solution as a new product line, Centera.

The benefits of an EFDMS are significant – by providing a true enterprise computing environment, data access times can be measured in seconds rather than days or even weeks. The reduced response time for the user community can accelerate the development cycle, additionally reducing overall costs. While reducing cost and increasing productivity, the skills of your staff can be further enhanced to assist in the advancement of defense programs. Staff members once dedicated to data distribution through media duplication can be redeployed to contribute to the overall goals of the team, improving team moral and performance, and further enhancing the process of converting data into information.