# Moving On: When It Is Time to Re-archive

*Michael Selway; Quantum Corporation; Middletown, MD/USA*

## Abstract

*Data archives are built to last a lifetime. More than a lifetime. They are in place to protect all the data that man and machine create. And yet, today, the archiving storage technology itself is rushing ever faster into obsolesce and every year it becomes more difficult to move this Archived information into the next generation of storage and support expanding requirements for accessibility. The evolution of storage media used to be measured in years, now it is measured in months. So how do you plan for migrating exponentially increasing amounts of data [1] when the subsequent data migration may have to begin even before the current migration is complete?*

*Planning a tiered data storage paradigm shift was never easy even when there was sufficient time to review options; to select the proper next step in what must be a never ending cycle of legacy data protection. But when time is short (and growing shorter), there must be migration strategies and methods established to provide not only for the next migration, but also setting the stage for all subsequent migrations. Considering strategies for migrating "like-technologies", such as disk-to-disk or server-to-server, can often be staggered to avoid a more massive disruptive migration.*

*But when a migration requires the core Data Management Archive Software (DMAS) technology be replaced, the impacts ripple across the total tiered data architecture with the potential to impact the user community relying on that data.*

*This paper explores aspects of strategies about how a Tiered Data Management Environment (TDME) facing the prospect of replacing the underlying DMAS should design an approach for converting and re-casting the archive.*

## Overview

Managing organizational data assets is an exponentially growing complex problem made more difficult by increasing organizational user (the "Community") demands to see and use that data. Shrinking life-spans of storage media complicated by increasing needs for greater storage make data migration soon an annual occurrence. With many moving parts in a TDME paradigm, some migrations are more readily accomplished than others. Methods and technologies exist to make actual storage migration relatively non-intrusive to the Community. Migrating across legacy to new disk drives or across legacy to new storage tapes can be planned and scheduled. This provides clear understanding of how and when the Community will be impacted. There exist numerous measurable parameters to assess and estimate the costs involved [2]. But when the underlying management software for organizing that data across disk and tape tiers must be replaced, the level of organizational disruption is high, likely unavoidable and ultimately very expensive in terms of personnel/IT resources consumed. Regardless of the approach taken, Data Center Managers (DCMs) are confronted with costs and disruptions to business as they move the

mountain of data assets. However, within the framework of this massive undertaking, efficiencies in approach and considerations of the overall system can significantly reduce the impact of the migration. This becomes the time to also examine the impacts on converting the form of the data storage as well as the access and handling of the data, to re-examine how the data is being used, if the rules and policies still match the intended protections and access, and decide on the core DMAS to underpin the TDME moving forward.

## Why is there a need to move on? What happened to the notion of "infinite" archiving?

The driving factor to migrate tape-based data from one format to another has not changed over the decades. Initially, the movement was prompted more often from a concern for the stability of the storage media [3]. But as technologies made all media sufficiently "reliable" (when compared to the expected useful life of the associated media drive technology), the impetus shifted to primarily reducing data management personnel expense, media costs, and sustaining environmental expenses for storing more information. Active data storage (on disk technology) requires migration to denser, faster technologies for the same reasons as Archive data storage (on tape technology). But migrating Active data is a less severe problem than migrating Archive data. Graphically, Active data may be viewed as growing at an increasing rate (as represented by single points along the *line* curve in Figure 1), however, the amount of archived data is growing at a far more rapid rate (as represented by the *area* under the curve in Figure 1).
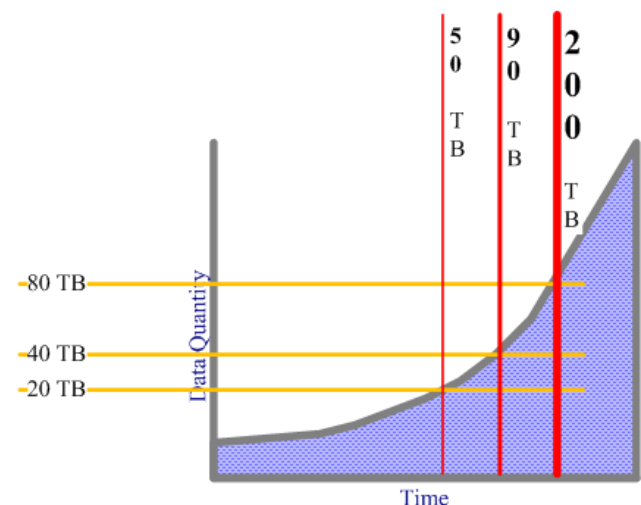


**Figure 1**

As a result, major disk vendors have means and approaches for handling disk media migration and these are well

documented and understood [4]. But few DCMs are eager to approach migrating archived tape data managed with DMAS to the next generation tape. There are four main reasons; lack of standard data-on-tape formats (dozens of tape drives and technologies exist today), lack of standard migration approaches, the shrinking time interval between required migrations, and the more complex level of indirection which exists between the archive data files and the file system that manages them.

migration across disks is also available from many disk vendors resulting in no disruption of data access to the Community. As the management of disks is also relatively less complicated, the DCM might choose a slow evolutionary migration of the more critical time/access-sensitive data to the newer disk architectures while continuing to retain the older disk suites for less critical active data. Later, as operational requirements dictate, the older disk drives will be decommissioned and removed from the architecture.
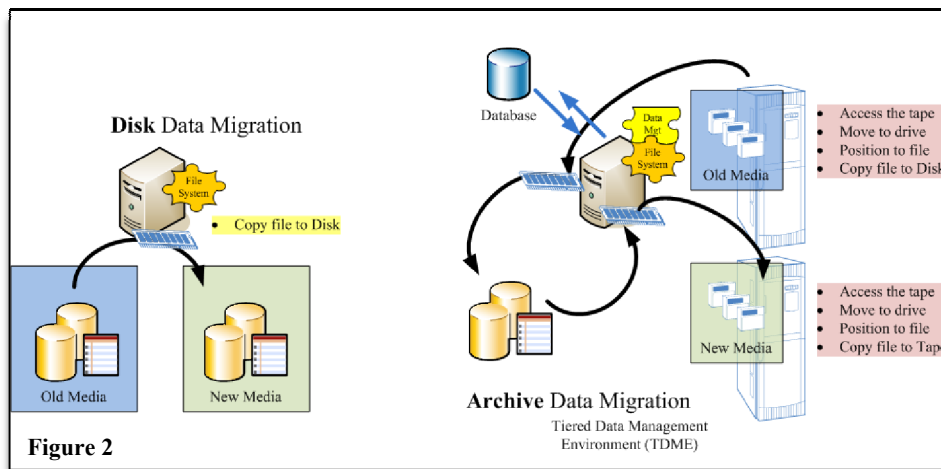


**Figure 2**

## *What to Move?*

As shown in Figure 2, disk-based data can be readily migrated to another disk-based system because the actual data resides directly under the file system managing it. But by the nature of tape-based archival data being managed as part of a comprehensive TDME the location of data is maintained virtually.

The linkage to data must be maintained in tables and data structures separately. This is the primary architecture of DMAS; to virtualize the apparent location of data such that applications/users always access it from the same "place" but the physical location is managed by the DMAS. Consequently, the simple copying of one tape to another does not result in the same location addressing therefore data access from above the file system would be lost. To maintain continuity, the data locations on the new tape must be passed to the TDME to update its records. The problem is further complicated when the migration includes the DMAS being replaced. Non-standard means of locating data elements on tapes placed there by different DMAS requires translating the storage scheme as well as moving to new tapes. These factors combined have led to the inescapable requirement to read all the data off the legacy formatted media and store it onto new media in the new DMAS format. However, as discussed below, the methods and timing of the full data retrieval can significantly impact the total cost and duration of the migration.

## *Migrating the Disk Storage*

As depicted in Figure 3, for disk technologies, migration is less disruptive therefore there's less impetus to do so preemptively. If more capacity is needed, more disk drives are purchased normally of the newer technologies. Transparent data

## *Migrating the Tape Storage*

When tape is simply a peripheral storage location for directed data storage, e.g., copying a file from disk to tape, the ability to transcribe to newer tape technologies is also simpler. Similar to disk migration, a simple tape transcription will normally be a sufficient approach for a technology refresh as the location of data on the tape is a managed by the tape medium and tape access software directly. The format of TAR [5] is an example of technology-independent reading and writing of data.

But when tape technologies become an extension of the active data environment, a data migration is far more complex therefore the organizational resistance to undertake a tape-based archive data migration is usually higher. However, this
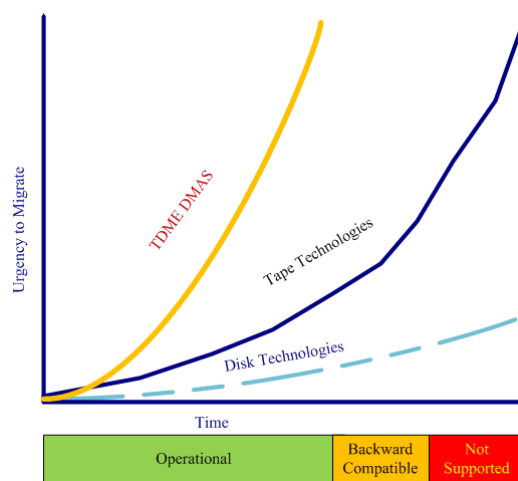


**Figure 3**

heightened resistance may be offset by a heightened urgency of concern for losing data access due to more rapid tape-drive technology obsolescence and more data requiring management. Unlike disk storage, the clock is ticking on how long data captured on legacy tape technology will remain accessible.

Tape drive vendors are competing to deliver higher capacity and throughput drives every 12-18 months. While most vendors will provide a measure of backward readability to their older formatted tapes within the same drive family (normally two technology revisions) the window of time to move off legacy tape formats before obsolescence is proportionally shrinking. This combined with the exponential increase in data being placed on tape (versus more linear storage requirements for disk) more data centers are crossing the point where the next data migration must begin before completing the previous one. This problem will not be easily resolved; ultimately DCMs either purchase more tape drives to do more migration stream in parallel faster, or constraints on the data being kept must be implemented. But this model is not sustainable indefinitely either and no new storage technology is yet widely available to replace tape-based archiving as the standard lowest cost per TB medium.

### Migrating the DMAS

As serious as tape migration is, an even more serious event that is often overlooked is when the underlying DMAS is determined to require replacement. This situation doesn't occur often and ideally would never happen. But when it does, the impact on the company is far more disruptive than a disk or tape media migration. The requirement to re-architect the TDME driven by a new DMAS structure will create data management disruption at the core level. To what degree this disruption is apparent to the Community varies significantly with the approach take for the migration. As normally the TDME is managed through a file system interface shifting the tiered DMAS may also necessitate working with a different Operating System and different file system. As part of the data migration this dictates creating a new file system and the data physically moved from the legacy to new file system. As the transition occurs this creates more administration challenges because the Community will need to work with potentially two separate systems in parallel.

Because of the broad impact on the organization and considerably longer migration planning cycle required when a DMAS replacement is involved, the CIO must monitor their DMAS Vendor for clues if/when the DMAS may be at risk of discontinuance. Waiting until the DMAS Vendor announces the DMAS is End of Life puts the organization behind schedule to migrate. To be preemptive, there are indicators, listed in Figure 4, that are the smoke before the fire.

But preemptive TDME migration also comes with added costs and creates disruption to a company's process and use of data because it impacts not only the archive data, but also the current active disk based data structures. Therefore, a DMAS migration effectively forces a disk-data migration plus a tape-data migration and a full re-examination of the TDME.



Warning Signs the DMAS Longevity is "At Risk"

- Parent company is purchased by a company not focused on Storage or Tiered Data Management
- Roadmap for future features not available or shows mostly vague general items
- Key new features previously planned for release are now withheld "indefinitely"
- Technical Support calls taking longer to get responses, even longer to get fixes
- Maintenance agreements shrinking to "one year" increments
- Maintenance plan options are being removed leaving only the most costly one(s) available
- Corporate literature and business reports no longer mention the DMAS or the data management solutions it provides
- Can no longer find Solution references on the Vendor webs site using the DMAS
- Key Engineering staff for the software development and technical support are leaving or being fired
- The parent company declares the DMAS to be in Sustaining Mode or other phrase translating to no new investment.

**Figure 4**

### How to Move?

When it is apparent that a majority of the TDME is requiring some level of migration to a new architecture, this also provides an optimal time to re-examine what the TDME is doing for the Community. In many cases, it has been years since such an analysis was undertaken. Community disruptions are going to occur at some level anyway; therefore making adjustments driven by refreshed understanding of the Community data access/storage requirements will only cause an incremental increase in total effort and costs. Further, based on what is determined, the total costs may actually be reduced due to developing new efficiencies based on understanding of what data is required to be archived and what data needs to be accessible today. Migrating to a new DMAS is also the time to re-examine how the data is being used and if the rules and policies still match the intended protections and access. Updates in policies can also lead to more efficient configuration of storage resources.

### What are the Requirements for the TDME today?

Organizations often become demand-driven in the areas of upgrading/enhancing their TDME. The original design met some or all of the requirements defined at the time. Budget for the physical (component infrastructure) and logical (storage management policies and media assignment) design were computed, allocated, and implemented. The system support structure of facilities, administrators, operators, allocation of compute and storage resources were defined parts of the Community and integrated into the TDME operations. And for a short period of time, the TDME met the stated needs of the Community. However, the Community is not static and almost immediately there are pressures to change.

The initial step therefore is to reexamine the expectations and users of the TDME. Interviews and surveys need to be conducted to determine what the system needs to provide for services and what information/data must be obtained and preserved to provide those services. Based on the results of building a new mapping of what the TDME must provide, updates in areas of the hardware components should be examined and added in as part of the TDME upgrade. In almost every case, new clients will be added to the system, new technology such as virtual machines used for server/client systems can replace discrete components, networks are faster and reach further, and total data management will evolve based on new Community expectations, SLAs, and government / industry mandates for security. Along with these changes come new upgrades for the data center itself. More efficient architectures will result in a change in facilities layouts, the heating/cooling/power design, and the number and mandated skill-level of the staff required for managing the new TDME.

### Picking the DMAS

Crafting a re-design of the TDME will still require the migration of existing data. But the rules and policies governing its use and protection may also guide the choice of the new DMAS and associated upgrades in disk and tape technologies. The selection of DMAS with underlying server and operating system should be primarily based on the capabilities required of the TDME by the Community. For example, Community data access tools and applications may be designed to work with the current DMAS structure. Therefore the replacement DMAS must either provide the identical appearance or the Community applications/tools will need to be adjusted. The latter situation will add time and expense to the migration effort.

But when choices are otherwise similar, the deciding factors should center on the educated expectation of the longevity and viability of the Vendor owning the DMAS. With infrastructure and storage technology changing rapidly, the CIO must be confident of the Vendor's commitment to support and expand the DMAS for at least the next five years.

### Budgeting the TDME Upgrades

The cost calculation and justification represent a separate analysis predicated on a clear understanding of what the current environment costs the organization (Total Cost of Ownership (TCO)). This value is crucial to showing the justification and true business value for upgrading the TDME. Combine those expenses, including the DMAS, which are viewed as new additions within a Capital Expense (CAPEX) plan. Add to this the expected Operational Expenses (OPEX) over the next five years. Importantly, compute and assign values to the improved organization of what can be done better, faster, less labor, supporting increasing new business efforts, etc. These are more qualitative characteristics, but to the organization their benefit values are at least as important as the infrastructure costs.

Finally, the estimated cost of conducting the actual migration needs to be calculated. Determining the time it will take to complete the migration and costs associated with the disruption to the Community can be calculated to a reasonable degree of accuracy. Figure 5 lists some of the more common

considerations, but every site will vary to how significant each factor may be. "Downtime" in particular should be careful examined as recent studies reflect the "hidden" cost to revenue measured in thousands to millions of dollars per hour [6].



Factors to Consider in Computing Data Migration Time and Expense

- Do the calculations for each media, drive type, and file system
  - Number of drives (total, available for migration effort, adding drives)
  - Files per file system (counts, average file size)
  - Performance for reading files (rated speed, adjusted speed for mounting, positioning, dismounting)(not count media compression)
  - Capacity of tape (total, estimated average actual in use)(not count media compression)
  - Total hours per day availability of resources to do migration (24/7? 5/8? Afterhours?)
  - Time to create a single file on the new DMAS (files per second)
- Aggregate the values for file systems and derive estimated number of days to complete
- Calculate total costs of migration accounting for
  - Operator/administrator hours and rates
  - Additional personnel resources during migration
  - Loss of productivity of the system to the Community (cost of running in Degraded Mode)
  - Loss of Revenue when data is unavailable (downtime for transitioning) [6]
  - Level of testing and evaluation of the transitioned data required
- Add in "factors for the unexpected" (delays, outages, unreadable tapes, etc.)

**Figure 5**

Still, this calculation of total costs cannot be an exact science therefore significant effort should be expended to ensure each factor is as accurate as possible.

Do not estimate conservatively on computing time and resources required to complete the data migration. Invariably, migrations take longer than expected due to unplanned factors occurring outside of the planned approach. In truth, any number of factors (such as tape drive/server outage, tapes being partially or fully unreadable, retrieved files are found to be damaged, unforeseen higher priority processing that delays migration) can extend this transition indefinitely.

However, combining these risk and transition factors realistically will show that maintaining the current TDME until a data migration is forcibly required will over the longer term be more costly to the organization than approaching the inevitable migration and TDME update proactively and preemptively.

## Approaches for moving the data

When undertaking a data migration, access to non-migrated data and migrated data must be managed for the full duration. This requires extensive management for conveying to the Community where their data is actually accessed from; the original or the upgraded part of the TDME. The result is effectively the overall TDME operating in a "degraded mode"; not operating at optimal performance or providing optimal services to the Community. Minimizing the migration time is crucial to returning the TDME to full capabilities but data can only be moved so fast. Therefore, a key aspect of a data

migration is determining means for minimizing disruption of data access to the Community.

Some approaches can do this more effectively than others. The following sections discuss three of the primary approaches for performing a DMAS migration. Note these approaches may be combined to varying degrees to create additional approaches.

### #1: Read it all back, write it all out

The most commonly considered approach is the traditional read the data from the legacy DMAS into the file system underpinning the new DMAS. Approached systematically, as depicted in Figure 6, files are identified into related groupings, read back over the network, and written in a parallel file system structure. This approach presents more planning difficulties however because throughout the migration Community files are split between two active DMAS systems. For a user looking for data, they need specific guidance as to where they should write new data and where to read the legacy data. Confusion and lost files can result from this approach. Working to assure a perfect migration, the DCM must perform comprehensive scheduling of what files to move and when they are to move.

As the data migration is occurring in real-time to the Community ongoing need to access the TDME, the migration is often scheduled for times where conflict is minimized. This may mean working late night or weekend shifts, or scheduling multiple maintenance windows denying access to the system. Very labor intensive and the productivity of the Community will suffer from the confusion.
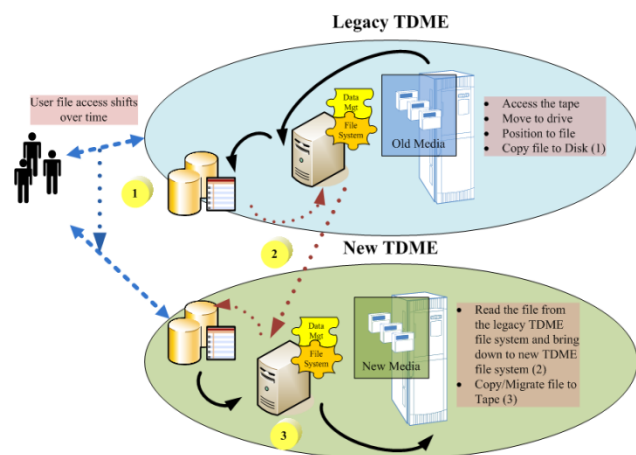


**Figure 6**

This approach will work well for a small environment of limited files and archive media requiring migration. But today most TDME's have large file counts, large numbers of archive media, and the environments are growing rapidly. As such, this simpler approach will not be feasible for organizations as the time to complete a large-scale migration iteration will eventually overlap the start of the next generational migration. Overlapping migrations significantly increase the confusion of users wanting to access their files, the planning for what data migrates, and increases the costs associated with the evolving TDME.

### #2: Split the TDME; maintain both legacy and new systems

A less organizationally confusing migration approach actually does no overt migration at all. Relying on the characteristics of the existing data, some organizational systems can adopt a passive migration approach. One scenario is the data being managed has a fixed life-span such that at the end, the data is deleted (or otherwise removed from the TDME control). A second scenario addresses the need to retain data for compliance purposes but without intent to actively use it. In these situations, there is not a compelling requirement to move the legacy data into the new TDME as doing so would also expend portions of the new TDME resources. Therefore, instead the legacy system is permitted to continue to exist in parallel to the new TDME.

As depicted in Figure 7, the division of data access is clearly demarcated and the Community understands that data after an established cutover time will be accessed from the new TDME. This approach also requires a file system structure that supports the organizational split of files associated with the legacy and new TDME's.
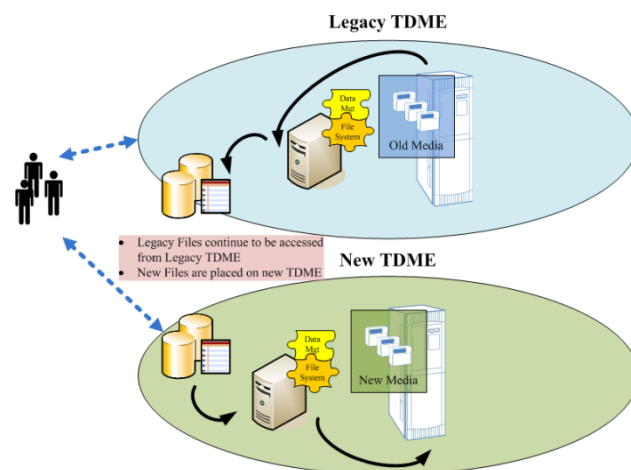


**Figure 7**

For scenario one, files in the legacy TDME will over time be deleted. Eventually, all the files will have been deleted at which time the legacy TDME and its DMAS can be decommissioned and resources re-purposed. In scenario two, the legacy TDME is expected to remain "active" (read-only) and maintained for an extended period of time. This adds costs to the total architecture to maintaining two discrete DMAS systems as well as added infrastructure and administrative overhead. A cost savings might be to re-purpose those components of the legacy TDME that do not directly support the read-only access paradigm. However, eventually, the ongoing operations of the legacy TDME technology will become a critical factor as the hardware vendors stop providing support. As a result, architectural migration will eventually still have to be performed if the legacy data is to remain viable in the organization.

### #3: Split Migration into 1) Namespace Conversion and 2) Data Migration

Select DMAS systems offer means to eliminate data access management by actually converting both legacy and new storage locations to a single access method. There are several strategies that can be constructed to do this, but the most effective one is replicating the full legacy file system namespace (the file metadata and the directory structures) into the new file system and directing all Community access to the new file system.

This approach recognizes that actual data content will exist on legacy and new formatted media, so includes extra information in the newly created structure indicating where the referenced data resides. This approach also minimizes the disruption to the Community as once the Conversion is done, access is the same regardless of where the data resides. Migration can now become a background effort occurring even during normal hours as the Community is effectively shielded from the data migration effort.

As depicted in Figure 8, the actual location of the file is transparent to the requesting application or user. As the storage methodologies of the legacy DMAS are typically quite different
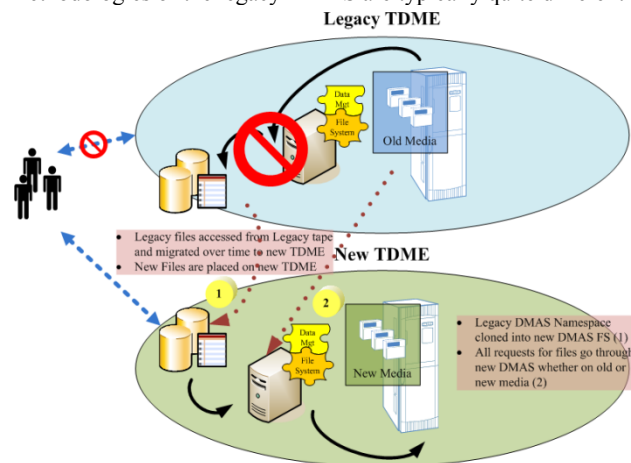


**Figure 8**

from the new DMAS approach, the ability to access the legacy format requires a deep knowledge of how to find the data on the legacy media. To facilitate the parallel access and minimally disrupt the standard operations of the new DMAS the most efficient approach requires a special conversion module. This module is designed to take "key identifying" metadata about a legacy file, determine where the file resides, and direct the mounting, positioning, and reading the file from the legacy formatted media.

This approach requires specialized knowledge to create the conversion module which most organizations do not possess. There are, however, technical companies that specialize in this type of migration and can offer partial or full conversion as well as data migration services.

Costs for the total migration are also reduced significantly because after the Conversion is completed, the legacy DMAS can be shutdown resulting in savings on maintenance and support. Later, after the entire legacy data has been read into the new TDME, the conversion module can also be removed.

## Summary

The decision to initiate a TDME migration is not an easy one and one that all DCM/CIOs would rather never have to do. However, technology changes and the evolution of business around DMAS vendors require this occur every 2-3 years. Unlike a more direct technology refresh of the disk or tape subsystems, re-examining the total TDME architecture driven by the underlying requirement to replace the supporting DMAS is invariably disruptive and expensive. The magnitude of disruption can be reduced to the Community (who are the most critical consideration) if the crucial *access* to their data is managed before actual *migration* of the data. Blending the legacy and new data access into the new DMAS structure establishes a "one place" view that is crucial to subsequently minimizing disruptions during the actual physical  migration.

When selecting the replacement DMAS focus first on what functional options each DMAS product provides as compared to the features needed for ongoing and expanding Community usage. From that list of qualified candidates, select a DMAS that affords the least disruption to the Community via a two-stage approach of consolidating the data access first, then undertaking the data migration over time. Through this 2-stage approach, the total migration costs will be less in terms of disruption and lost productivity of the Community.

## References

[1] The Diverse and Exploding Digital Universe, John F. Gantz, Project Director, *IDC*, March 2008. See http://www.emc.com/collateral/analyst-reports/diverse-exploding-digital-universe.pdf

[2] Reducing Costs and Risks for Data Migrations, Patrick Allaire, Justin Augat, Joe Jose and David Merrill, *Hitachi Data Systems*, February 2010. See http://www.hds.com/assets/pdf/white-paper-reducing-costs-and-risks-for-data-migrations.pdf

[3] Storage Media Life Expectancies, John Van Bogart, *National Media Laboratory*, June 22, 1998. See http://nssdc.gsfc.nasa.gov/nost/isoas/dads/presentations/VanBogart/VanBogart.ppt

[4] Example: See http://ams2300.com/nl/products/storage-software/multiplatform-backup.html

[5] TAR command, Linux Information Project (LINFO), July 15, 2006. See http://www.linfo.org/tar.html

[6] Understanding Downtime, Business Continuity Solution Series™ , *Vision Solutions*, 2004. See http://www.dsscorp.com/prod/dsscms.nsf/m/Vision%20Solutions%20Understanding%20Downtime.pdf/$file/Vision%20Solutions%20Understanding%20Downtime.pdf

## Author Biography

*Michael Selway, Senior Storage Software Consultant at Quantum Corporation, has worked with tiered data storage management systems since 1984 where he was part of the creative team that architected the E-System / EMASS FileServ HSM software. Mike went on to spend 14 years designing and consulting on tiered data management systems centered over LSC SAMFS. He continues this role with Quantum StorNext /Storage Manager shared data management software, designing innovative and visionary tiered file management architectures.*