

# The Role of Active Archive for Petascale (and Beyond)

Addison Snell

November 2010

*White paper*

## MARKET DYNAMICS

Legends get amplified in the retelling. Each generation passes its lore to the next, and the peril may grow with each successive iteration. The jungle grows deeper, the cliff gets steeper, the teeth become sharper. It seems it just takes more drama to impress an audience these days.

One story the supercomputing audience can seem bored of hearing is the need for data management structures that can keep pace with advancements in microprocessors. In this case there is reason to be increasingly dramatic. Data sizes continue to grow at exponential rates that outstrip even theoretical computational improvements, and meanwhile the requirements for storing and accessing data for longer periods of time are growing as well.

In the high performance computing (HPC) industry, storage spending continues to increase at a faster rate than server spending<sup>1</sup>, despite persistent reductions in storage costs and improvements in data density. We've been talking about the data management monster for decades, and in this case, tales of its ferocity are not merely hyperbole. This ogre really has grown larger, hairier, and more terrifying every year.

Intersect360 Research studies indicate that HPC sites on average are experiencing 40% annual growth in both the number of files and the total amount of data stored. There are compound reasons for this trend, including:

- *Requirements for more accurate models* – Both scientific and engineering design models generally require ever greater fidelity in order to advance research and to design next generation products.
- *Increase in model complexity* – Model complexity is driven both by requirements to include more components (e.g. larger molecules in computation chemistry, large sections of an automobile in crash analysis), and involve more complex analysis approaches.
- *Expanded and improved data sources* – New data is being provided from sensors, new and improved instrumentation, greater numbers of measurement devices, more complex experiments, etc.
- *Increase in data-oriented research* – Researchers are using data filtering and correlation techniques to identify interesting results from new data streams and to re-examine existing data. These sorts of applications can generate intermediate results that can add to the total volume of data.

As the world's top supercomputers progress into petascale computing, they will continue to encounter more treacherous hazards impeding the path to productivity and insight, including achieving adequate bandwidth, reliability, and density for increasing levels of scale. The intrepid pioneers will continue to find innovative ways to tame the data management beast. In the petascale era, one such technology will be active archival.

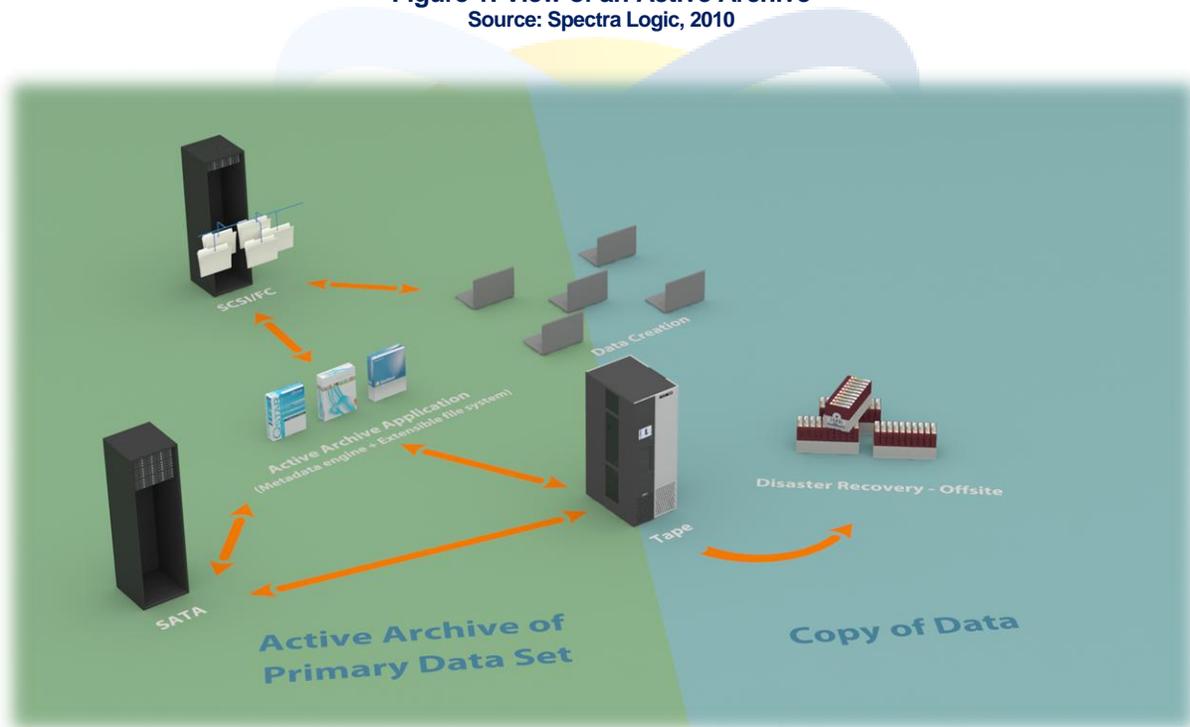
---

<sup>1</sup> Intersect360 Research HPC Market Advisory Service, "Traditional HPC Total Market Model and Forecast: 2010 to 2014," June 2010.

Tape is part of many large-scale HPC storage infrastructures, and it will continue to be. For the largest data centers, this is a simple matter of practicality. Tape tends to be significantly less expensive than disk when measured per terabyte, per square foot, or per watt consumed. Optimizations in the storage deployment can free up more capital budget and more physical plant resources for other HPC technology elements, such as server nodes and interconnects.

Active archives reinforce the economic value of tape to high-capacity, high-performance storage environments. An active archive contains the complete metadata indexes and data searchability of an extensible file system, improving visibility to data on either tape or disk, and optimizing access to it. [See Figure 1.] This helps prevent data that is usually “WORN” (write once, read never) from becoming “lost in the archive,” in the sense that you’ve forgotten you have it or don’t know how to access it. HPC data tends to represent the core intellectual property of an organization, and active archives help organizations continue to leverage this knowledge.

**Figure 1: View of an Active Archive**  
Source: Spectra Logic, 2010



To promote the awareness and adoption of active archives, Spectra Logic, a leading supplier of tape archives, joined with other industry members to form the Active Archive Alliance. This multi-vendor group is working to establish file system, application software, and hardware standards for building and utilizing active archives.

An effective active archive architecture can help execute information lifecycle management (ILM) schemes, migrating data as appropriate from primary disk to more cost-effective archival, yet still retrieving unexpectedly necessary pieces of data as needed. Consider the following examples:

- In-use data, such as that used in hurricane tracking, must be continuously updated and available, requiring access times that may be milliseconds. Store this data on disk either online or near-line.
- Video-on-demand servers must be able to display everything that is available in an organized, sortable fashion, and can tolerate access times of several seconds. Store this data near-line.
- Long-term, large-scale, archived data, such as seismic data from a two-year-old oil exploration effort, can tolerate retrieval times of a few minutes. Store this data on tape, which can be accessed in one to two minutes.<sup>2</sup>

These examples imply characteristics that an active archive must incorporate to be effective for petascale systems. First of all they must be able to efficiently store large volumes of data, whether that is driven by fewer, larger files or numerous smaller ones. Efficiency in this case includes the consideration of both the capital costs (e.g., continuously lowering dollars per byte) and the facilities costs of the infrastructure (e.g. increasing density and optimizing power consumption and cooling).

Second, the storage system must be able to handle large data sets coming into the archive without significant performance bottlenecks. As data sizes scale, throughput becomes an important consideration, and for many supercomputing-class applications it is essential.

Third, the archive must be “future-proofed” for longer time periods than other parts of the data center, because the nature of the problem is that it must hold increasing volumes of data for significant periods of time. This includes the cost to maintain the components over time, as well as compatibility to future technology generations. By incorporating both disk and tape in an active archive, organizations can mitigate the risk of technology obsolescence.

## INTERSECT360 RESEARCH ANALYSIS

In the petascale era and beyond, there is an even greater need to plan and deploy technologies for the management of large-scale data. Intersect360 Research believes active archives will be one of these critical technologies. By improving visibility and access to data, active archives help organizations manage large volumes of data over time, potentially reclaiming datacenter resources to go toward computation.

Active archive technology is delivering a new level of capability to large-scale computing environments, allowing them to store and easily retrieve data in ways that were not previously possible. Incorporating tape into an HPC environment no longer requires a sacrifice in data access or searchability. Active archives empower scientific and technical computing users to continue to use tape to manage the perpetual growth in data. As a founding member of the Active Archive Alliance, Spectra will continue to help enable the productive and critical role of tape in supporting scientific advancement, as more HPC users leverage their historical knowledge more efficiently and future-proof their new findings for the next generation of discovery.

---

<sup>2</sup> Spectra Logic has demonstrated access times of 65-75 seconds.