

---

## Collaborative Data Life- Cycle Management - Until We Invent the Storage Time Machine

---

*Arun Jagatheesan*  
(*San Diego Supercomputer Center*)

### *Abstract*

A “Storage Time Machine” is a hypothetical storage media that can provide infinite storage space (capacity), extreme I/O benchmarks (performance) and relatively low cost to manage large-scale data over the years (cost of ownership). We know such a storage media with all those features does not exist. Until we invent the “storage time machine”, data storage infrastructures will have to tackle the I/O discrepancy problem, and also the need to manage data stored in multiple storage media with various I/O bandwidth rates. This problem is exacerbated in collaborative data lifecycle management environments, where multiple distributed organizations with heterogeneous storage technologies participate in a collaborative infrastructure to manage data lifecycle. This is a growing trend in large-scale scientific projects and global enterprises that operate in multiple continents. The LSST project is expected to manage 200+ petabytes of replicated data, distributed in several countries. Assuming an ideal world with relatively low network latency, a HPC I/O in any data center could theoretically be routed to another data center in a different continent. In addition, multiple replicas of the data could provide different scenarios on how the HPC I/O could be achieved. How do we effectively use multiple heterogeneous storage resources to serve broad communities that want to share petabytes of data collections? In this talk, we discuss the relevance and requirements of HPC I/O in collaborative data management. In addition, we discuss research topics such as policy-based-I/O and the need for new data protocols for large scale end-to-end HPC I/O.

---