**CRAY**

# LANL and Cray Answer a Performance Problem with DataWarp Accelerator Innovation



Los Alamos National Laboratory (LANL) has a weighty mandate — solve national security challenges through scientific excellence. Chief among them is to ensure the safety, security and reliability of the U.S. nuclear stockpile as part of the National Nuclear Security Administration's (NNSA) Stockpile Stewardship program.

This responsibility drives yet another challenge — meeting current mission needs with appropriate compute technology while also adapting to the larger evolutionary and revolutionary technology changes that will shape future simulation environments.

In response, NNSA's Advanced Simulation and Computing program established an initiative to develop and deploy a series of Advanced Technology (AT) systems. Their first installation is Trinity, a supercomputer based on the Cray® XC™ series architecture.

Using Trinity, LANL is exploring compute technology so they can provide platforms with higher performance for predictive capability.

**Challenge**
High-resolution 3D simulations generate extremely large datasets. Used to ensure resiliency during long-running simulations and for data analysis and visualization, these datasets are valuable but were limiting application performance and impacting throughput.

> "This gives us an order of magnitude increase in I/O performance with extreme predictability. We've never had this level of performance predictability before."
>
> *- Galen Shipman, Computer Science Lead,*
> *Eulerian Applications, Los Alamos National Laboratory*

The application runs for extended periods of time — often several months on end. As a result, data flushing to the Lustre® file system creates inefficiencies because the application stops completely during the flushing process.

Restrictive data flow and slow processing speeds (I/O bottleneck) made matters worse. In order to address these issues and enable their applications to run as efficiently as possible, LANL had only two options: 1) frequent check-pointing with reduced recovery periods but highly inefficient run times; or 2) bare minimum number of check-pointing runs and long recovery periods.

LANL had reached the financially feasible limits of scaling with traditional HDD technology configurations and needed a better solution.

## Solution

LANL approached Cray to build a solution based on flash storage. Their goal was to achieve speeds of 3.2 TB/s, enable hourly flushing to the Lustre file system and shorten recovery periods.

Cray proposed an idea Gary Grider, division leader of the HPC division at LANL, had begun pursuing several years prior. That idea was to develop an SSD storage product consisting of service nodes connected directly to the Aries™ network, each containing two SSDs and an API/library with functions to initiate stage in/stage out and query stage state. The solution could be configured in multiple modes using the workload manager.

The resulting solution is the Cray® DataWarp™ applications I/O accelerator.

### SYSTEM DETAILS

Cray® XC™ series supercomputer

41.5 PF peak performance

2+ PB memory capacity

19,420 compute nodes

78 PB parallel file system

3.7 PB burst buffer storage

The DataWarp accelerator is unique in two ways: 1) it uses the Cray Linux® Environment and DataWarp software runs in root level; and 2) applications can call a library in real time, thereby circumventing a secondary scheduler.

## Results

With DataWarp technology, LANL has achieved:

- 3 TB/s goal
- 15x improvement in I/O performance
- Reduction in checkpoint restarts to 60 seconds (down from tens of minutes)

> "With DataWarp we have the capability to efficiently index and analyze truly large datasets. We're finally unlocking the insight stored in data and enabling new science."
>
> *- Bradley Settlemyer, Senior Storage Researcher, Los Alamos National Laboratory*

The improved resilience and performance of large-scale, 3D, high-fidelity simulations has resulted in increased detail for visualization and analysis output. Furthermore, with better utilization of the system's computing resources the team can generate more comprehensive data output from their codes.

LANL is seeing "pretty staggering" efficiency at 1,024 nodes (32,768 cores). It means they have more bandwidth available. And more bandwidth means they can generate more valuable datasets with negligible overhead to the application.

"This gives us an order of magnitude increase in I/O performance with extreme predictability," says Galen Shipman, computer science lead on the Eulerian applications project. "We've never had this level of performance predictability before."