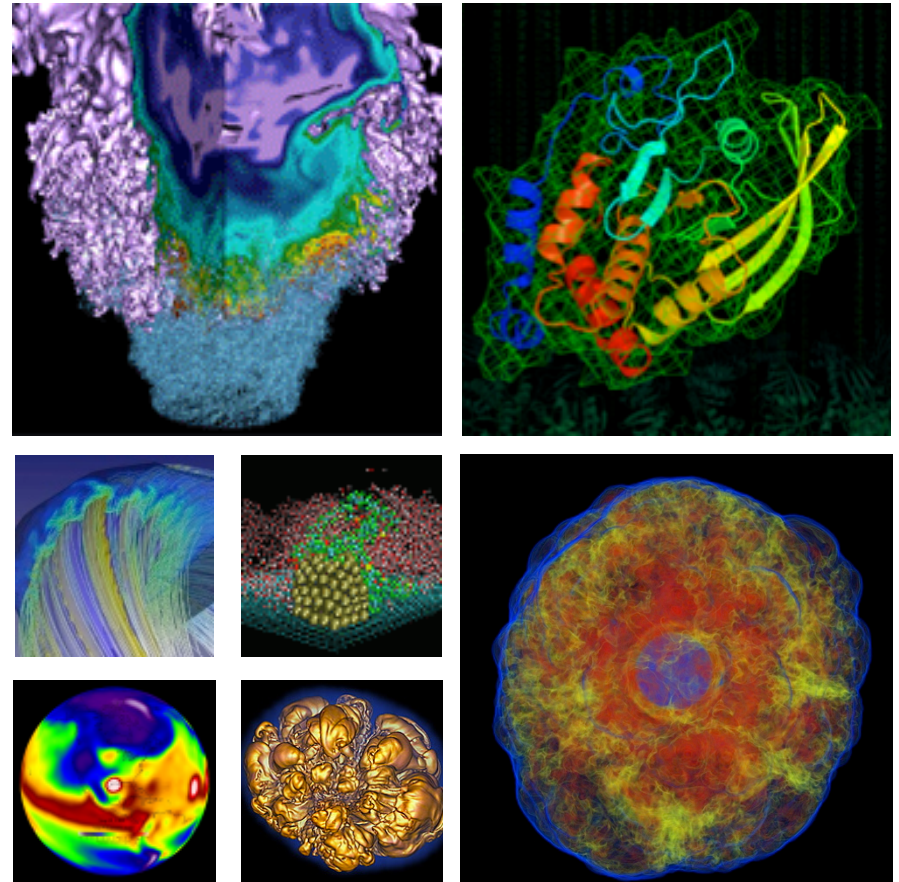


# Scalability Challenges in Large-Scale Tape Environments



**Jason Hick**

**Lawrence Berkeley National Laboratory  
NERSC Storage Systems Group**

IEEE MSST  
June 4, 2014

# Agenda



- **NERSC and its storage systems**
- **The Golden Age of Tape**
- **Our challenges at scale**
  - Reading data, system usability
  - Proactively maintaining the system
  - Having enough people
- **Industry challenges at scale**
  - Component and end-system reliability
    - Mechanical failures – flash, disk, tape
  - Speed versus size of single devices
  - Detecting and repairing failures
- **Summary**

# National Energy Research Scientific Computing Center (NERSC)



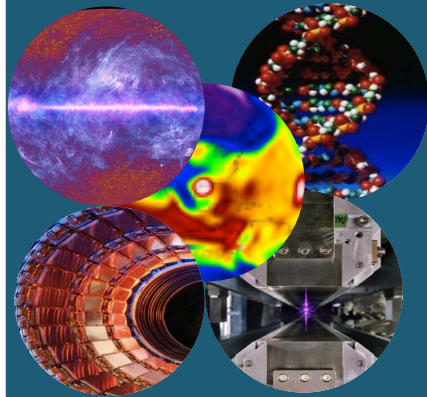
© 2014 The Regents of the University of California, Lawrence Berkeley National Laboratory

- Located at Berkeley Lab
- User facility supports 6 DOE Offices of Science:
  - 5000 users, 600 research projects
  - 48 states; 65% from universities
  - Hundreds of users each day
  - ~1500 publications per year
  - With services for consulting, data analysis and more

# Types of computing at NERSC

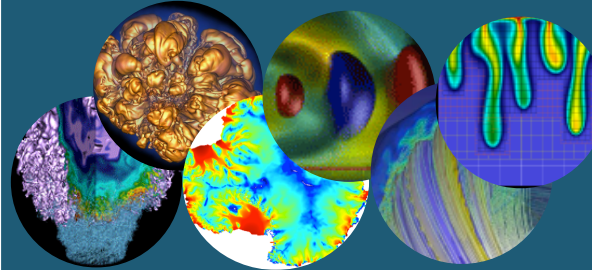
## Data Intensive

Experiments and Simulations



*NERSC ingests, stores and analyzes data from Telescopes, Sequencers, Light sources, Particle Accelerators (LHC), Microscopes, and other scientific instruments*

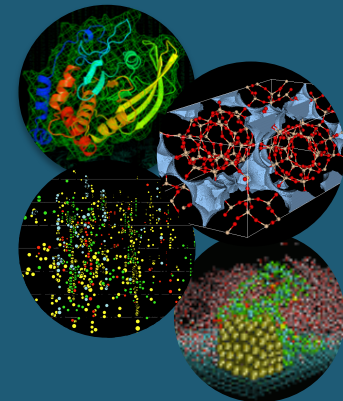
## Large Scale Capability Simulations



*Petascale systems run simulations in Physics, Chemistry, Biology, Materials, Environment and Energy at NERSC*

## High Volume

Job Throughput



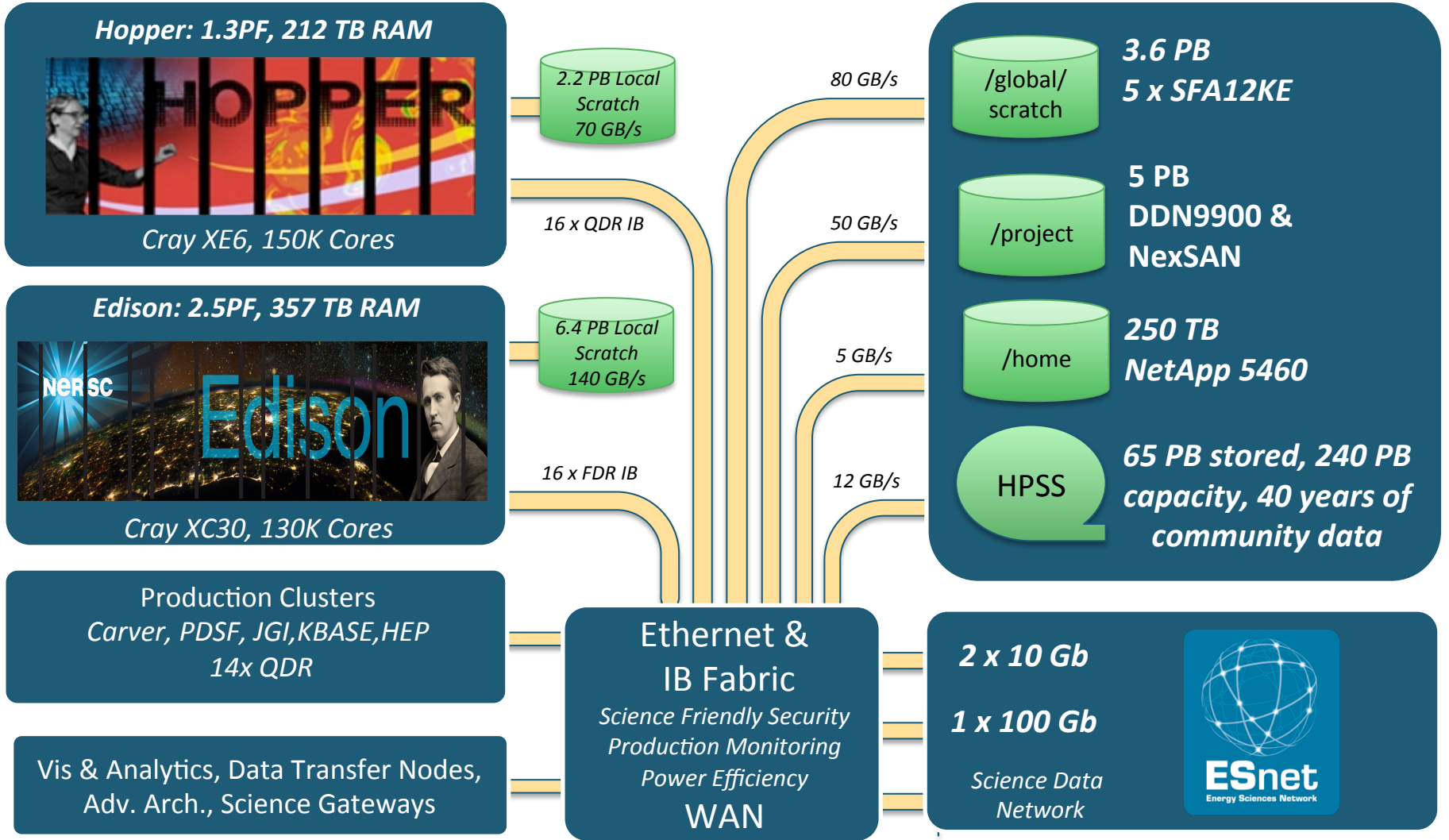
*NERSC computer, storage and web systems support complex workflows that run thousands of simulations to screen materials, proteins, structures and more; the results are shared with academics and industry through a web interface*

## NERSC

Petascale Computing, Petabyte Storage, and Expert Scientific Consulting



# The compute and storage systems 2014



# Focusing on storage at our facility



- **Parallel file systems (Lustre and GPFS) are primary storage to supercomputers**
  - Total of over 20 PBs of disk available to users
  - Some multi-PB parallel file systems backed up to HPSS (Parallel Incremental Backup System)
    - Has demonstrated it can process over 150TBs of backup data in a single day currently using direct-to-tape with 12 T10KC tape drives
    - On average, we complete a restore for a user about every other week
- **Archival and backup systems (HPSS) are secondary storage for users**
  - 65 PBs of data stored, growing at >1PB per month
  - 30% of user IOs are read/retrieve requests from archival storage, so a very active archive
  - Focus on reliability of the system for user data by:
    - Deploying solutions to proactively monitor and maintain health of user data, and environmental parameters necessary for tape
    - Actively migrating/moving data within the system

# The Golden Age of Tape



- **Tape is demonstrating capability for future decades (150TB), still handling vast quantities of data, and is integrated to varying degrees with file systems**
- **Storing data at scale (>5PB) tape is power efficient, fault resilient, and cost effective**
- **Today, “Tape” or “tape systems” normally means HSMs**
  - HPSS, pDMF, and SAMFS
  - Few workloads go direct-to-tape (e.g. large-scale backups, instrument/raw data acquisition)
- **Tape continues to enable the highest data growth rates**
  - Supports >50% CAGR at our facility
  - Address the most difficult data ingest, migration, and long-term storage needs

# Our challenge – reading data from tape



- **Unordered requests result in**
  - Wasted time mounting tapes
    - Especially in the case of large amounts of data or ingests over time where large quantities of tapes are involved
    - ~1 minute per tape mount
  - Longer duration of overall transfer
    - Repositioning within tape
    - Re-mounts of the same tape
  - Mechanical issues during excessive tape mounting cause further delay
    - Cartridge, tape library, drives
    - Mounts succeed but sometimes after multiple attempts
- **More important to store data to tape optimizing for your reads**
  - Even storing it in different ways
- **Things that were ok to do yesterday/year become a problem as time goes on**
  - Small files with tapes getting larger
  - Number of files stored over period of time (multiple tapes)

*Avoiding one tape mount is equivalent to reading a 12GB file.*

*Taking the time to order the list of files to retrieve can reduce duration of transfer from days to hours.*



# Our challenge - proactive maintenance



- **More devices and more complexity at scale**
- **Require advanced features and automation in problem detection, determination, and notification**
  - The industry only recently has software that helps determine problems with tape drive or media
  - We achieve, but struggle to be proactive
  - Proactively failing a tape reduces the duration and complexity of problem resolution
- **Vendors are moving away from onsite support**
- **Ability to detect and fix failures increasing slower than the system's capacity**
  - Validating a single tape
  - Rebuilding a tape or copying data off a tape

# Our challenge – having enough people

---



- **There is no metric on how many staff per PB, however more staff is required for higher complexity or scale of a system**
  - This is not unique to tape systems
  - Finding skilled staff is difficult
- **The mechanical nature of storage (disk and tape) makes them people intensive**

# Industry challenge – component reliability



- **Improve tape's environmental sensitivity**
  - Libraries are typically not sealed/filtered
  - Tape and drives are exposed to temperature, humidity, particulate in the room
  - Results in special considerations for tape, which gets costly
- **With capacity and quantity of tape increasing, need improvements in component reliability**
- **Failures are typically mechanical (cartridge, drive, robot)**
  - Reliability of data in practice on tape is very high
  - We observe that tape failures are not catastrophic as opposed to disk/flash failures

# Industry challenge – speed vs. size



- **Capacity leadership is vital to tape industry**
- **Need to balance tape drive speed with capacity improvements**
  - How many is too many files to risk putting on one tape?
  - How long will it take to rebuild the tape?
  - How long will it take to migrate the data off the old technology onto the new one?
  - How long will it take to verify a tape?

# Industry challenge – detecting & repairing failures



- **Software to improve on managing historical information/statistics**
  - Tape systems tend to be in place for long periods of time and have a wealth of statistics to understand
  - Move beyond break-fix into proactive health determination
  - Identifying suspect tapes (soft errors), drives with multiple failures, etc.
- **Tape drives are still configured by-hand upon replacement**

- **Tape technology is capable of enabling big data and exascale storage**
  - Highest supported CAGR, capacity demos show promise
- **Challenges at scale can be met with close collaboration/partnership of high scale sites and industry**
  - Improve software to detect, diagnose, and repair faults
  - Work to improve component reliability
  - Features to support tape ordering for high volume reads



**Thank you.**