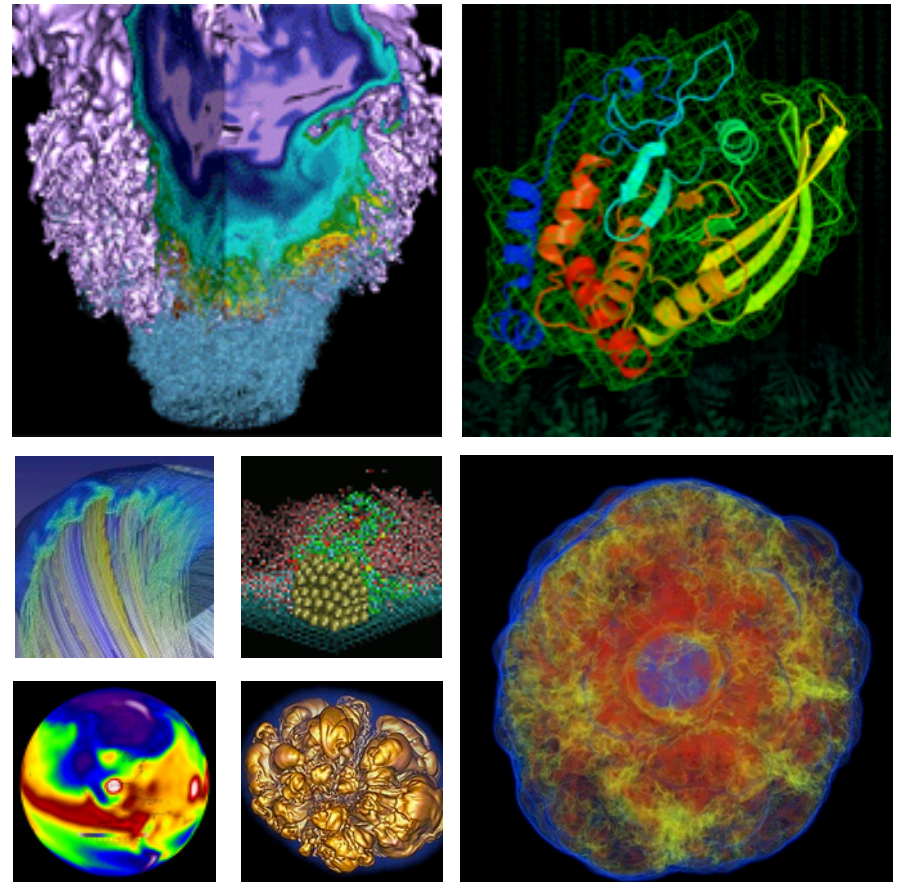


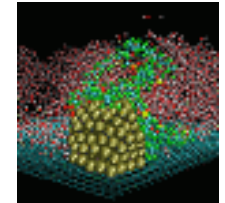
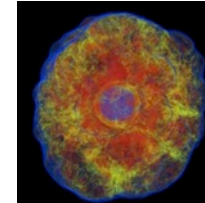
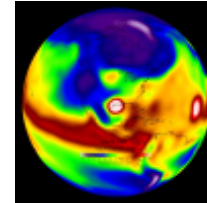
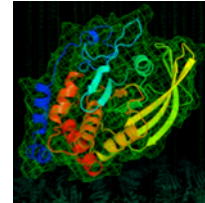
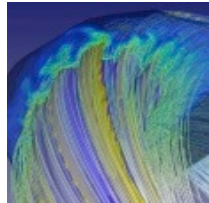
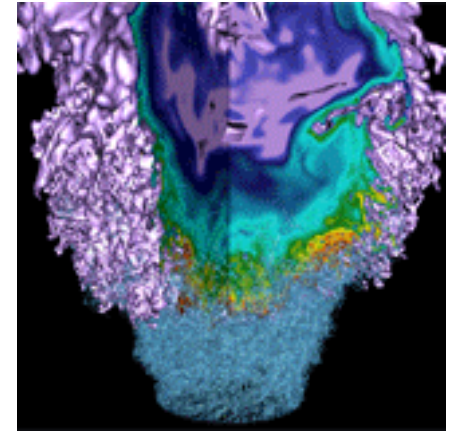
Data-Driven Science at NERSC



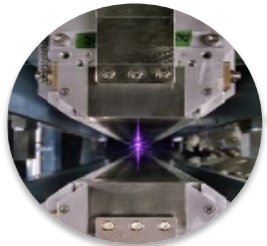
Richard Gerber
Senior Science Advisor to the Director
NERSC

August 8, 2013

New Community Focus on Data

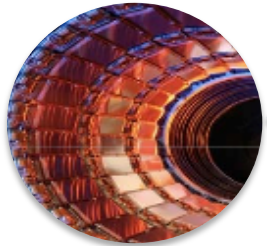


DOE has extreme data needs



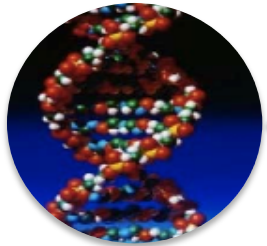
Light Sources

- Many detectors on Moore's Law curve
- Data volumes rendering previous operational models obsolete



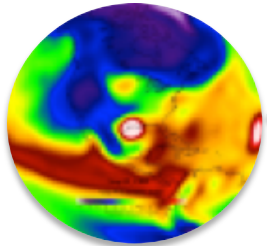
High Energy Physics

- LHC Experiments produce and distribute petabytes of data/year
- Astro: peak data rates increase 3-5x over 5 years, TB of data per night



Genomics

- Sequencer data volume increasing 12x over next 3 years
- Sequencer cost decreasing by 10 over same period



Computing

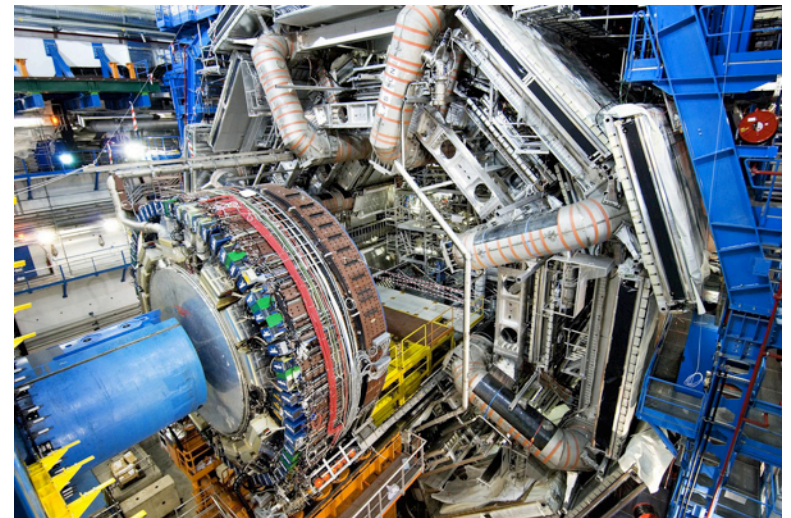
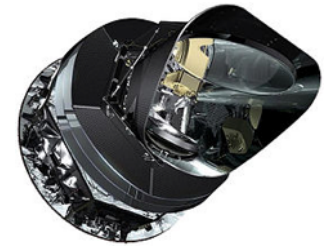
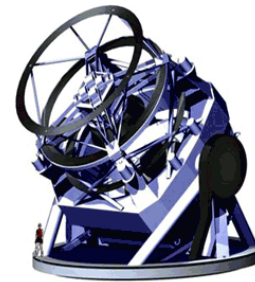
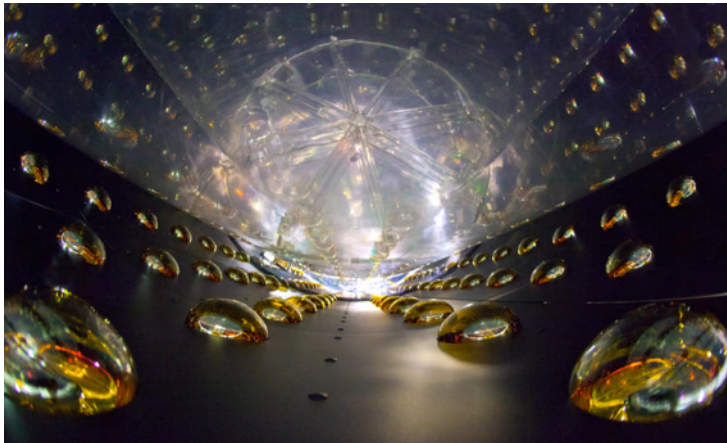
- Simulations at Scale and at High Volume already produce Petabytes of data and datasets will grow to Exabytes by the end of the decade
- Significant challenges in data management, analysis and networks

Source: Dan Hitchcock

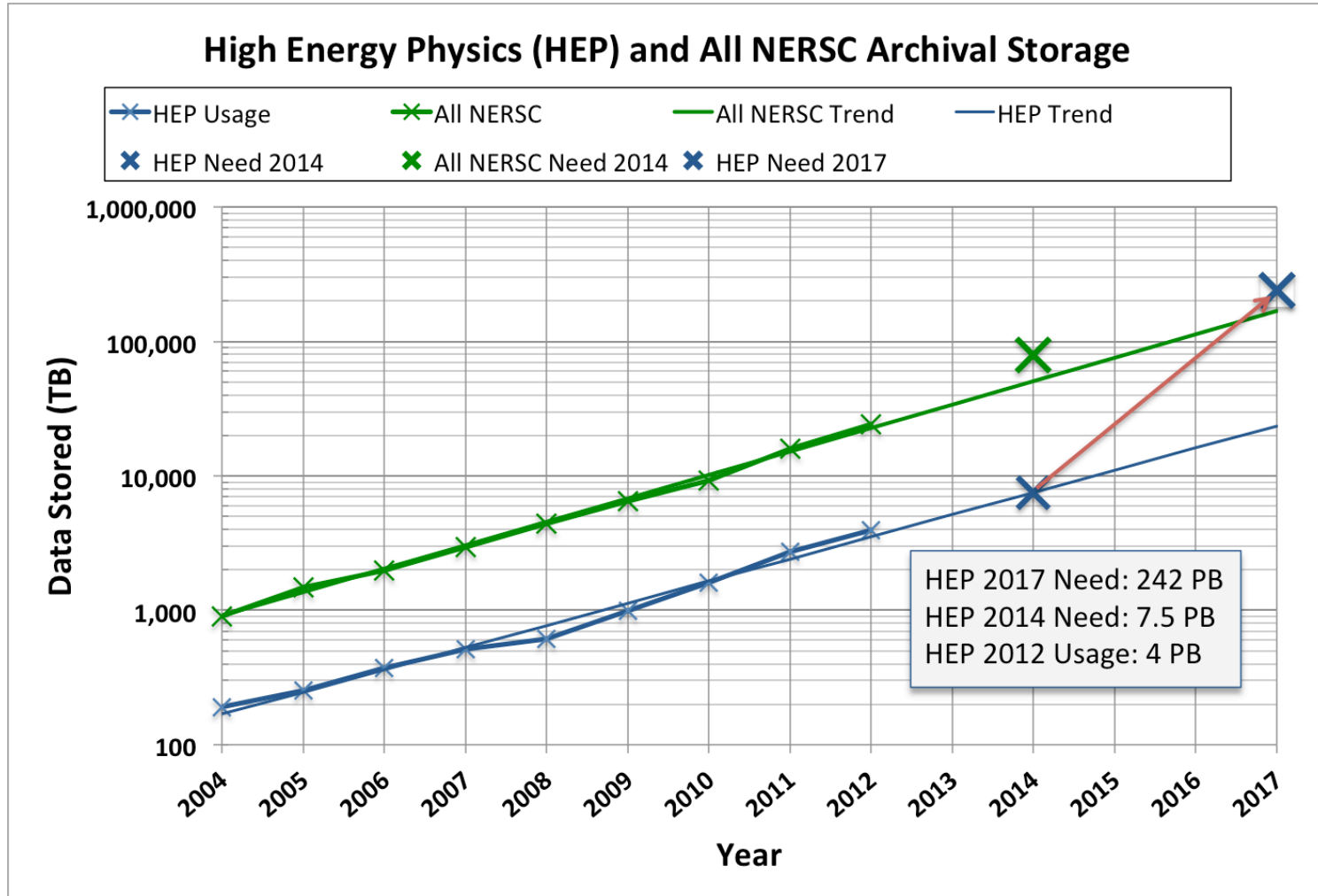
DOE Investment in Data-Producing Facilities

NERSC

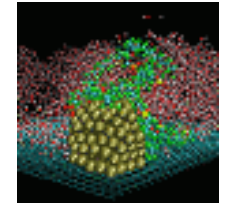
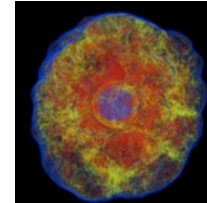
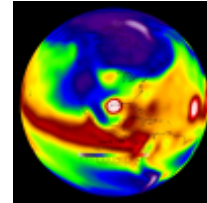
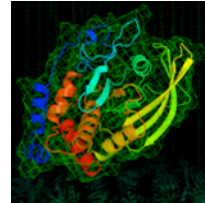
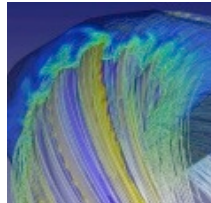
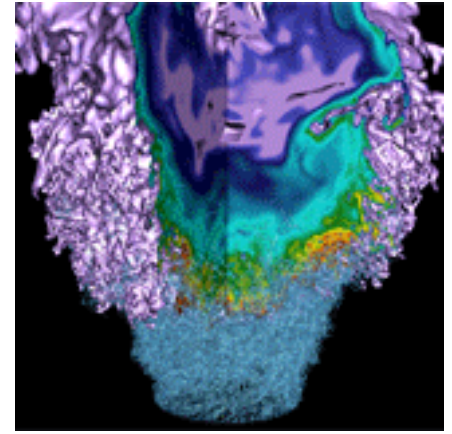
- DOE has a huge investment in facilities that will produce 100s of PBs of data



Exploding Data Storage Need Just in HEP

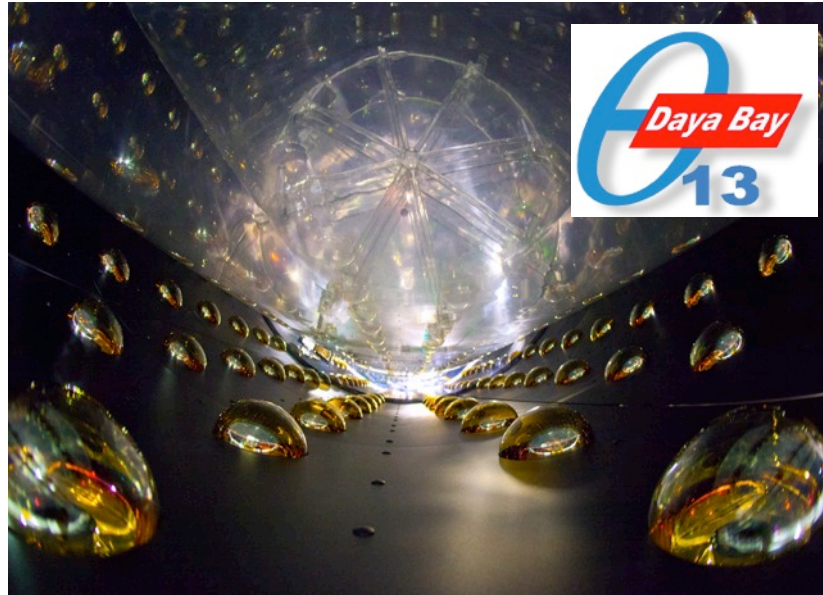


We're Already Doing a Lot



Discovery of θ_{13} weak mixing angle

- The last and most elusive piece of a longstanding puzzle: How can neutrinos appear to vanish as they travel?
- The answer – a new, large type of neutrino oscillation
 - Affords new understanding of fundamental physics
 - May help solve the riddle of matter-antimatter asymmetry in the universe.



Detectors count antineutrinos near the Daya Bay nuclear reactor in China. By calculating how many would be seen if there were no oscillation and comparing to measurements, a 6.0% rate deficit provides clear evidence of the new transformation.

Experiment Could Not Have Been Done Without NERSC and ESNET

- PDSF for simulation and analysis
- HPSS for archiving and ingesting data
- ESNET for data transfer into NERSC
- NERSC Global File System & Science Gateways for distributing results
- NERSC is the *only* US site where all raw, simulated, and derived data are analyzed and archived

The “*Supernova of a Generation*”



- Using NERSC facilities, including the Deep Sky **NERSC Science Gateway**, scientists were able to discover a supernova just hours after it first appeared
- Quick detection permitted immediate follow up observations, allowing scientists to gain unprecedented insight; SN 11KLY may eventually be most studied supernova ever
- Made possible by **state-of-the art computational and analysis tools, workflow, and a science gateway interfaces**
- Observations taken at Palomar Observatory (PTF) and data automatically transferred via ESnet



The supernova was discovered by Peter Nugent, a NERSC staff member, who leads the automated supernova search project.

- Funded and used by High Energy Physics, Nuclear Physics
- Networked distributed commodity Linux cluster in continuous operations since 1996
- **Detector simulation & data analysis**
- **Data intensive, high throughput workflows**
- **Grid Support**
 - **OSG, WLCG stacks**
 - **Compute and storage elements for OSG, ALICE**
 - **Storage elements for ATLAS**



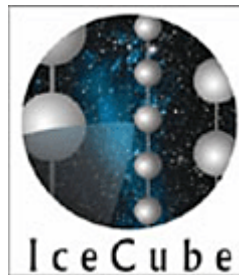
PDSF Quick Facts

- 2300 cores
- 1 PB globally accessible disk
- Interconnect 1GigE, 10 GigE, IB
- SGE batch system

PDSF Projects



- PDSF is an essential resource for a number of groups such as STAR (Tier 1), KamLAND, ATLAS (Tier 3), ALICE (Tier 2), DayaBay (Tier 1), IceCube, CUORE, etc.
- Groups such as SNO, SNFactory, CDF, BaBar, LUMI, Planck have used PDSF in the past.



- **Genepool**
 - Heterogeneous Mixture of “standard memory” and large memory nodes
 - Newest hardware is Intel Sandy-Bridge, 128 GB per node
 - Total of ~8500 cores (776 nodes)
- **Storage**
 - Nearly 5 PBs of storage
 - Mixture of legacy NFS and GPFS
 - Growing use of HPSS



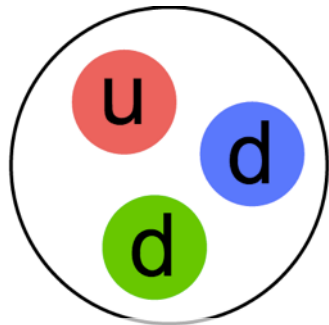
NERSC Data-Driven Science Gateways



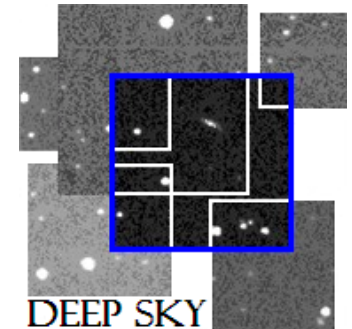
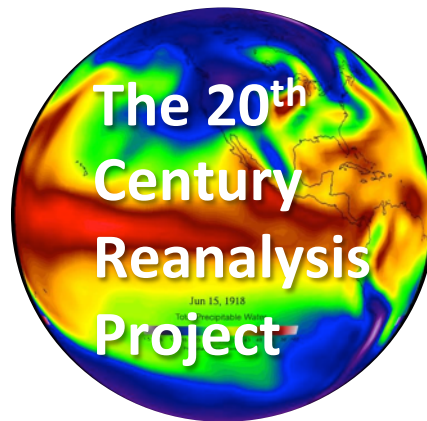
Simulation-Driven Data Science

Empirically Driven Data Science

MATERIALS PROJECT

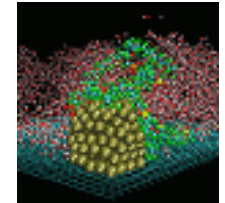
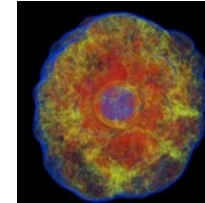
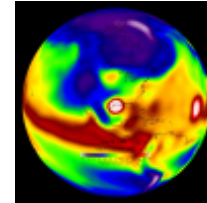
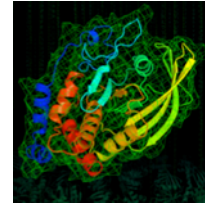
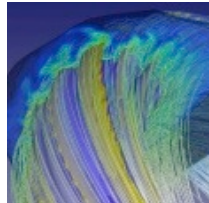
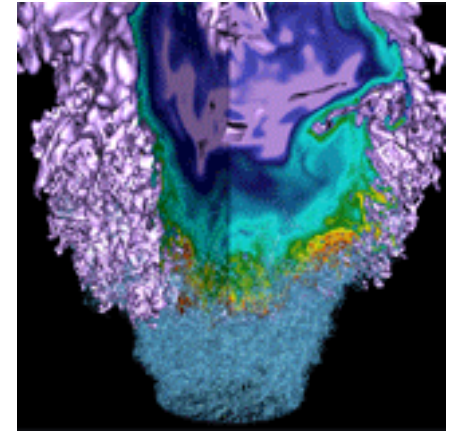


The Gauge Connection



OpenMSI: Advanced visualization, analysis and management of mass spectrometry imaging

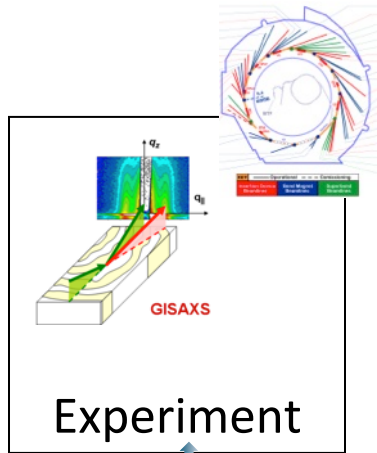
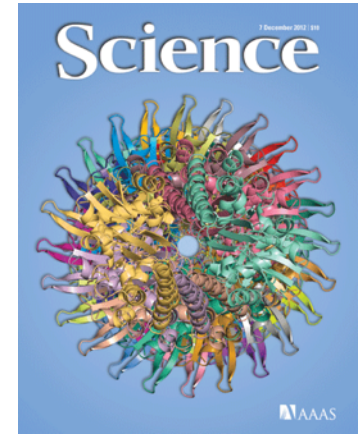
Where We'd Like to Go



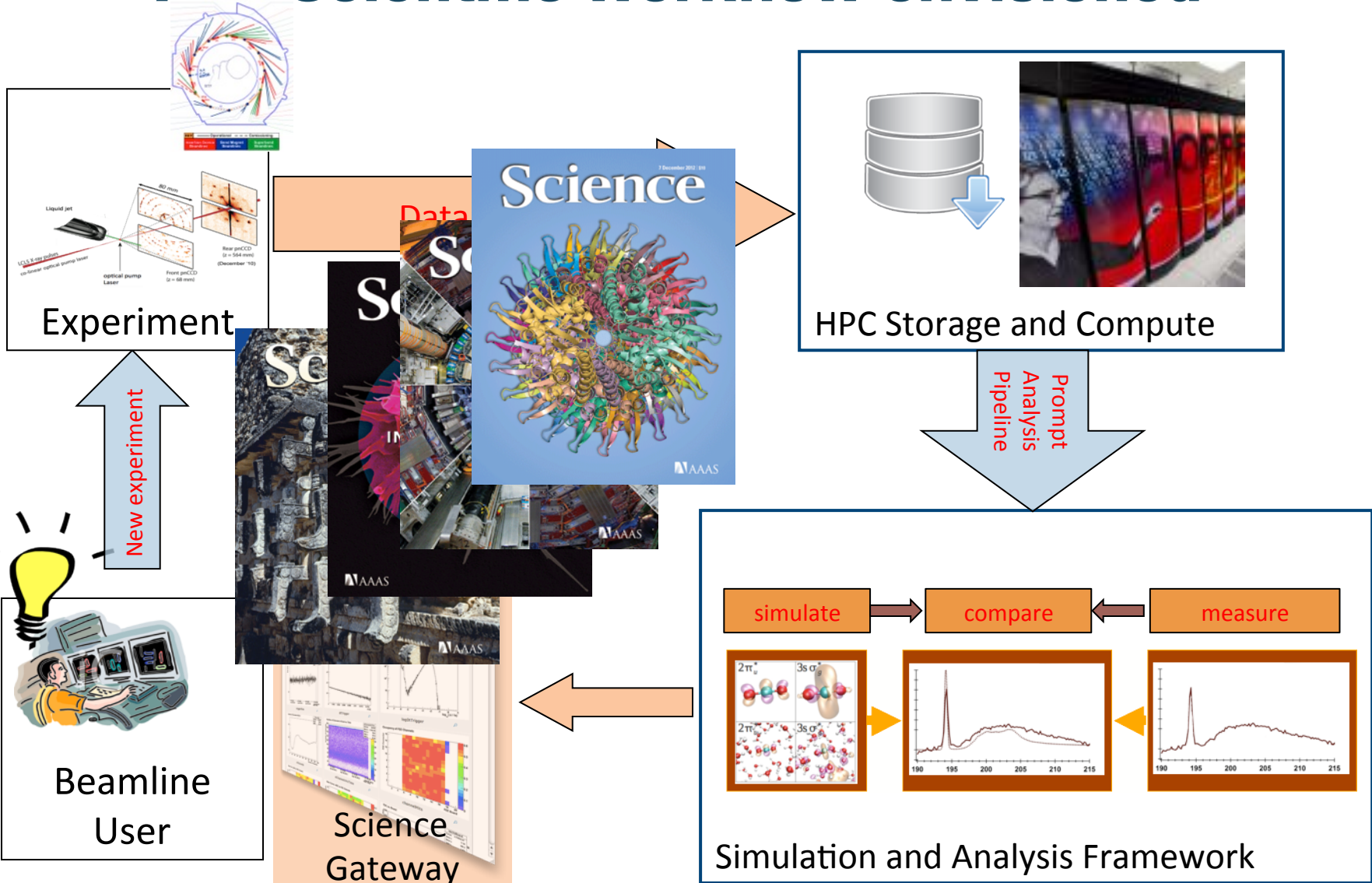
NERSC's Mission is to accelerate scientific discovery at the DOE Office of Science through high performance computing and data analysis.

Our challenge and our goal: enable scientific discovery through data by providing facilities and services not otherwise available.

ALS Scientific Workflow today



ALS Scientific Workflow envisioned



Our Challenge



- **Continue to support the data and I/O needs of the simulation community**
- **Enable scientific discovery through analysis of data**
- **Integrate simulation and analysis workflows**
- **Enable sharing of data among diverse communities using various technologies (web, shared disk spaces)**
- **Accommodate high-throughput workflows**
- **Provide training, consulting, and tools**
- **Do all this in a time of tight budgets**

- **Develop and deploy new data resources and capabilities**
 - Accelerate NERSC's traditional storage growth rate to meet rapidly increasing requirements for capacity and bandwidth.
 - expand our existing science gateways and data transfer systems to facilitate high-volume data intake and publishing
 - implement remote, redundant tape copies to provide increased protection for valuable data
 - deploy automated data management techniques to efficiently utilize storage
 - We are proposing to enhance the data processing capabilities of Edison in 2014 by adding large memory visualization/analysis nodes, adding a flash-based burst buffer or node local storage, **and deploying a throughput partition for fast turnaround of jobs.**
- **Partner with DOE experimental facilities and projects to identify requirements and create early success**
 - NERSC pilot projects have shown great success with automated data pipelining, indexing, searching, archiving, sharing, and distributing end users via the web

- **Provide expertise and services for extreme data**
 - Provide expert consulting and training
 - Embed “data postdocs” with science teams to develop new capabilities and train the next generation of HPC scientists (based on NERSC’s successful “petascale postdoc” program)
 - Develop sophisticated web-based gateways to interact with and leverage data
 - Support database-driven workflows and storage
 - Use scalable structured and unstructured object stores
 - Provide search and analysis software for massive data
 - Provide comprehensive scientific data curation
- **Leverage ESnet and ASCR research to create end-to-end solutions**
 - Deploy high speed networking between our data resources and DOE Facilities in collaboration with ESnet.
 - Use networks in software-adaptable ways to greatly improve data access, making it quickly available where it will have the highest impact.



National Energy Research Scientific Computing Center