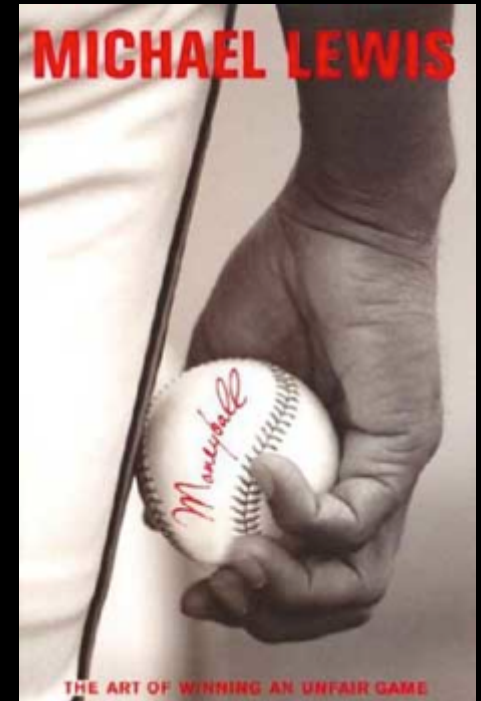# More Data, More Science, and … Moore's Law?

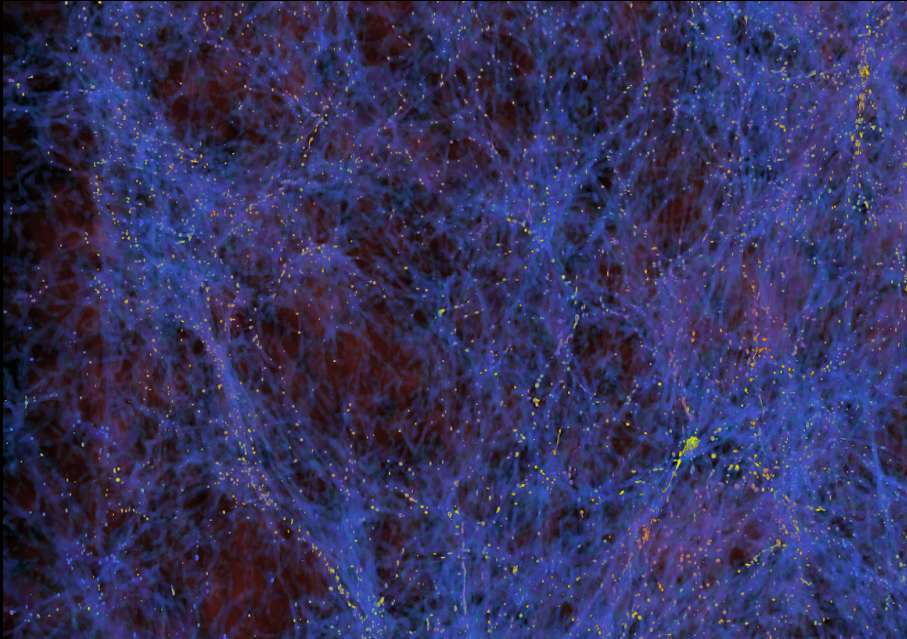## Kathy Yelick

**Associate Laboratory Director for Computing Sciences**
**Lawrence Berkeley National Laboratory**
**Professor of Electrical Engineering and Computer Sciences**
**University of California at Berkeley**

# "Big Data" Changes Everything…What about Science?

# Combine simulation and observation for next Cosmology breakthrough
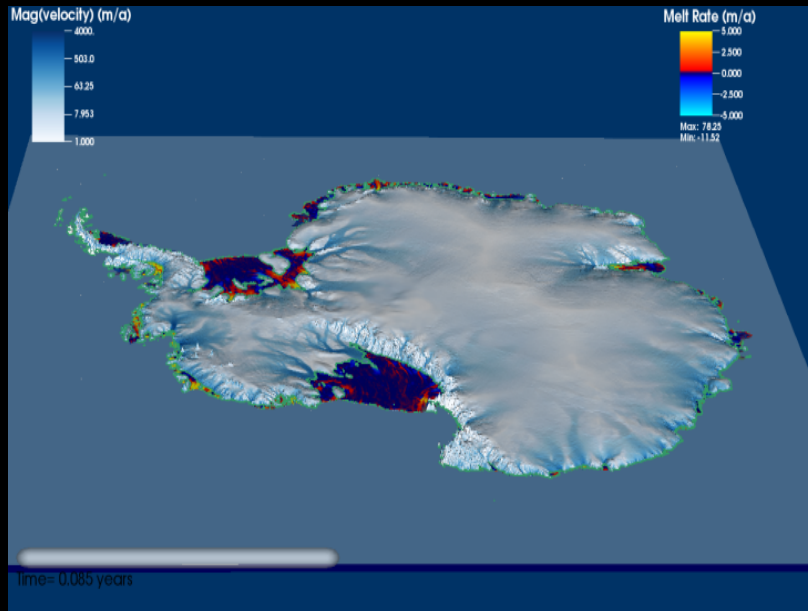


*Nyx simulation of Lyman alpha forest using AMR*



*Kitt Peak National Observatory's Mayall 4-meter telescope, planned site of the DESI experiment*

Reduce systematic bias in observation through simulation of ~1 Gigaparsec Baryon Acoustic Oscillations in the Lyman Alpha Forest and ~100 Gigaparsec simulation of galaxy clusters, both requiring adaptive mesh refinement (AMR).

# Climate models and microbial analysis together to predict the future of the environment



New climate modeling methods, including AMR "Dycore" produce new understanding of ice
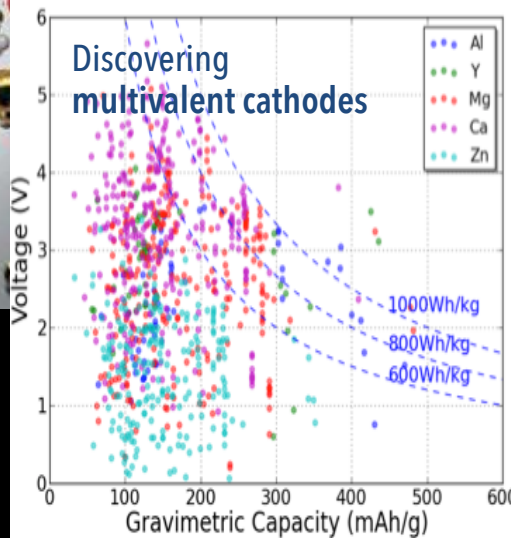


Genomes to watersheds Scientific Focus Area

Understand interactions between environmental microbiomes and climate change with *kilometer resolution models* that track dynamic 3D features (with AMR) and *genome-enabled analysis* of environmental sensors.

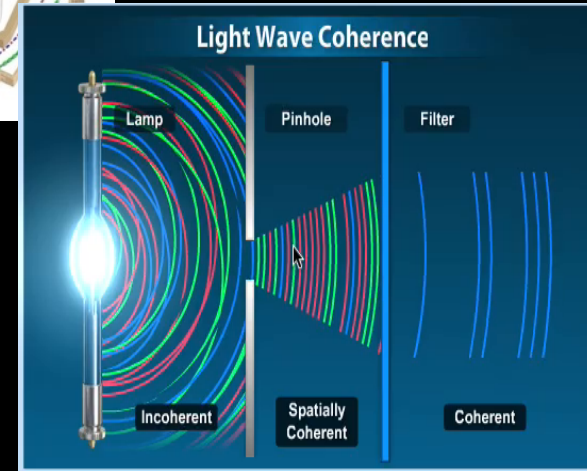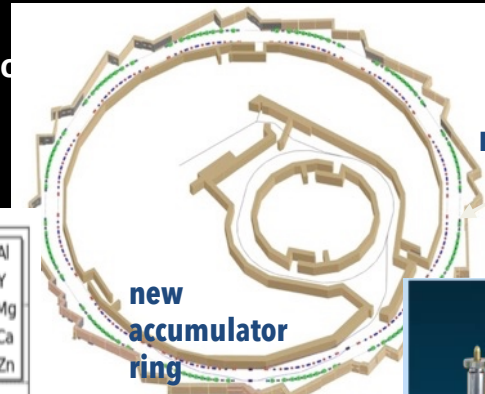# Understand and control energy with advanced light sources and materials modeling



**Materials Project**

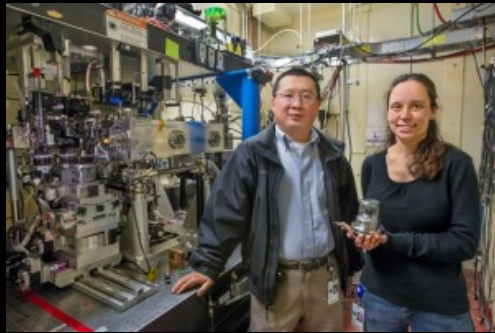**13,030 users hosted at NERSC with software co developed by CRD**

Discovering multivalent cathodes

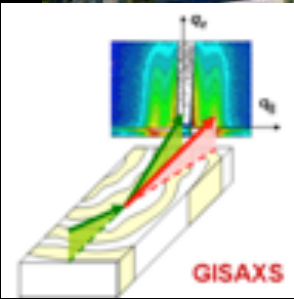ALS-U Upgrade

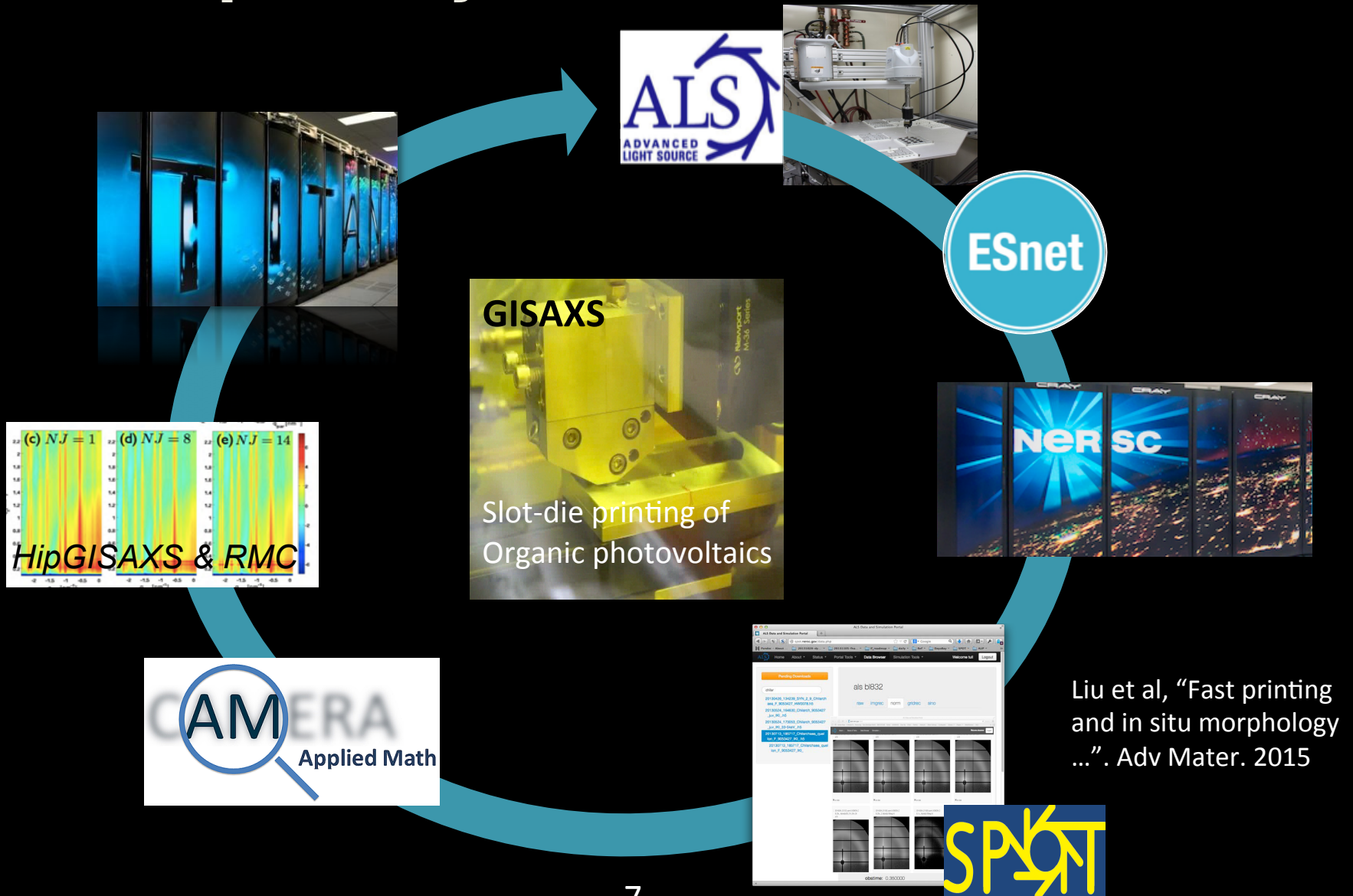new ALS ring

new accumulator ring

Light Wave Coherence

Understand and control the direction and flow of energy with minimal losses using *advanced instruments*, *high fidelity models*, and high throughput simulation and analysis for applications in energy, environment and computing,
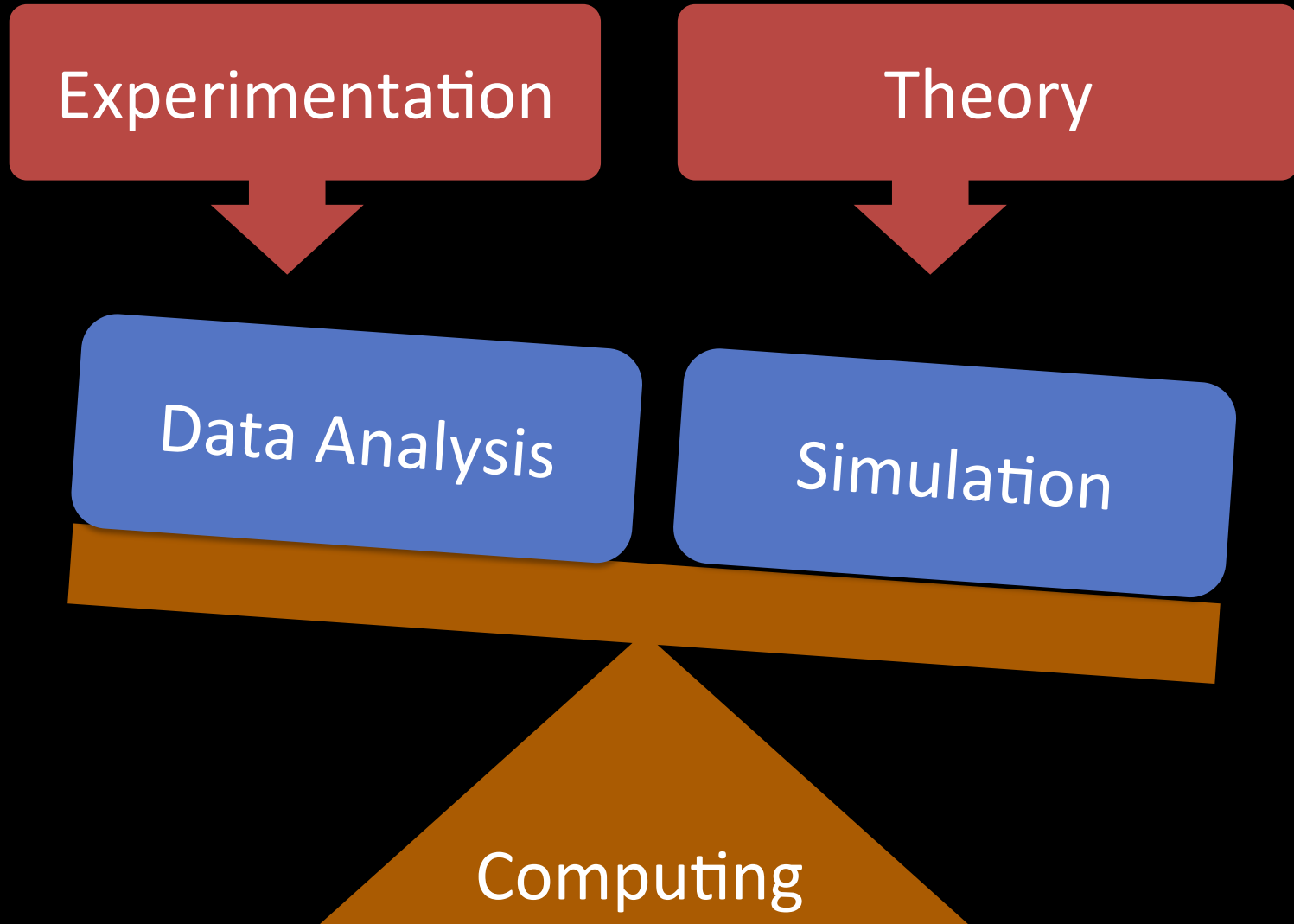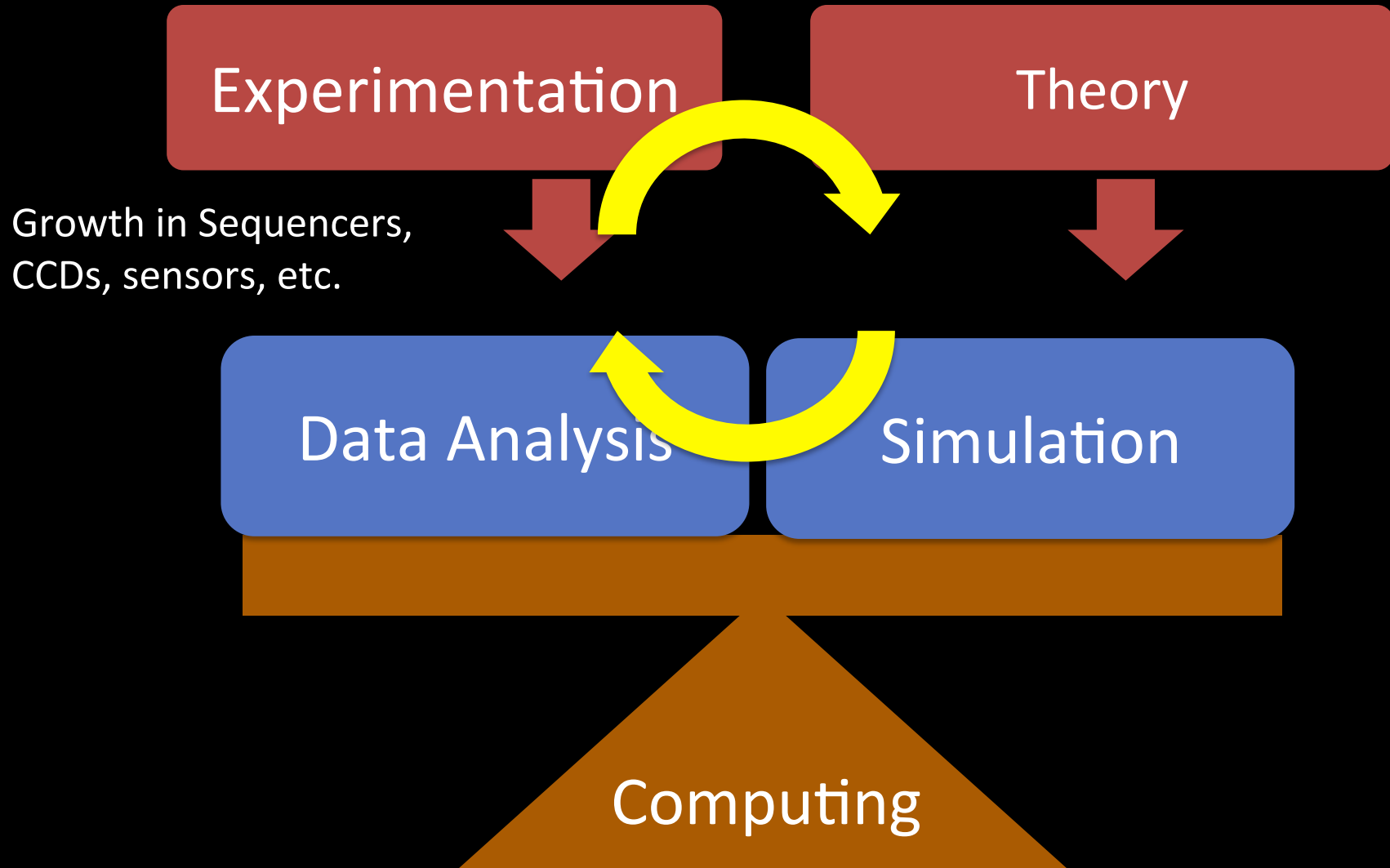
# Old School Scientific Workflow

# Computing, experiments, networking and expertise in a "Superfacility" for Science



**GISAXS**

Slot-die printing of Organic photovoltaics

*HipGISAXS & RMC*

CAMERA
**Applied Math**

Liu et al, "Fast printing and in situ morphology …". Adv Mater. 2015

# Old School HPC: only for Simulation

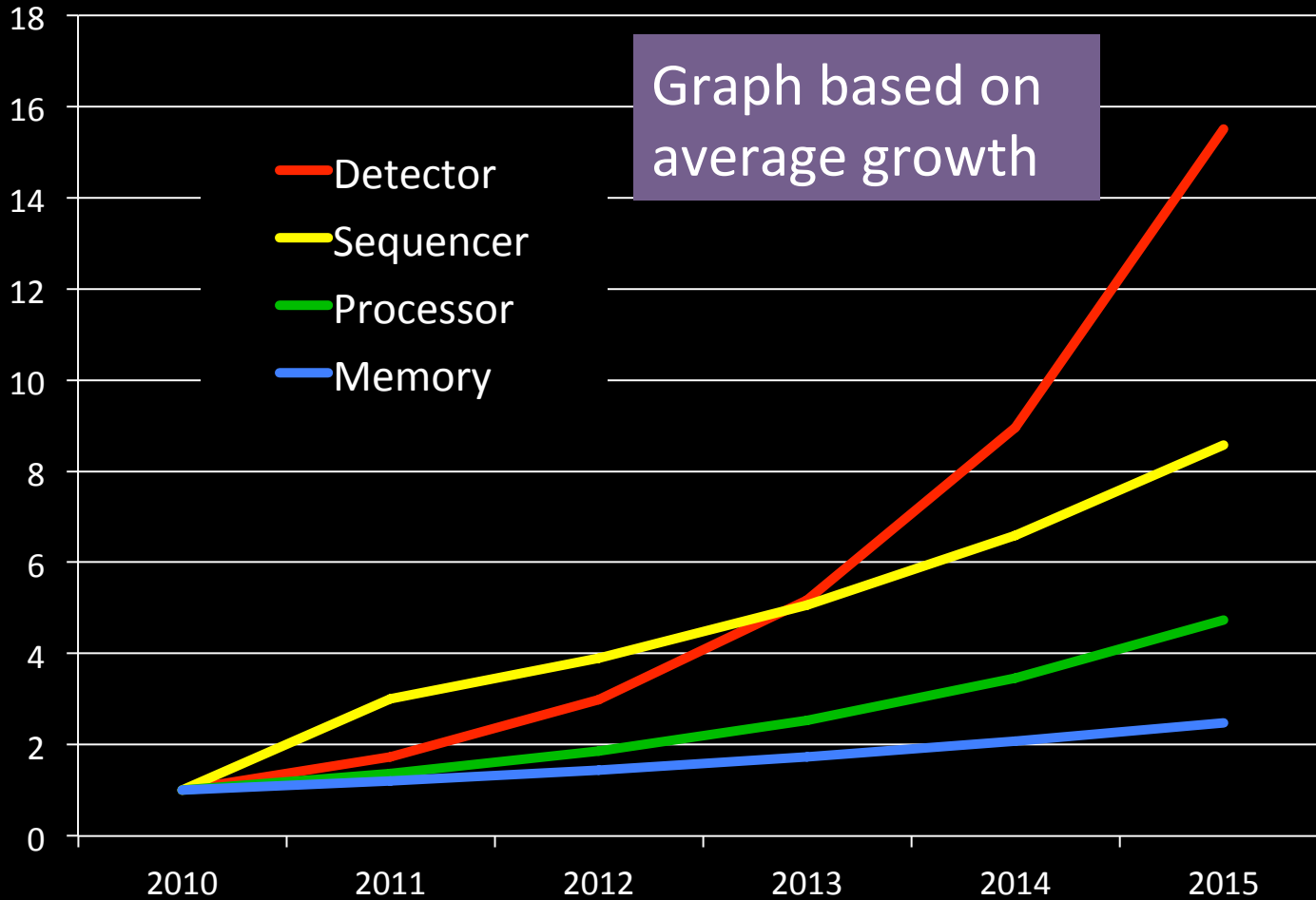# HPC is equally important in experimentation

# Questions?

1. **Are there MSU examples of "science at the boundary" of simulation and observation?**
   – How should you take advantage of these opportunities?

**Part 2**

# The Data Tsunami

# Science Data Growth is Outpacing Computing

# Old School Scientific Data Search

# Automated Search, Meta-Data Analysis, and On-Demand Simulation



Automated metadata extraction using machine learning

Jobs submitted by "bots" based on queries; algorithms extract informatics for design

# Questions?

2.  **What are the largest and most complex sources of research data at MSU?**
    – What types of data/CS/math/stat challenges are there?

**Part 3**

# Networking and Computing Facilities  Need to Adapt

# ESnet: Exponential growth in networking

Petabytes/month

Legend:
- Traditional IP
- Transatlantic
- Big science data

Y-axis: 0, 10, 20, 30, 40, 50, 60

X-axis: Jan-15, Feb-15, Mar-15, Apr-15, May-15, Jun-15, Jul-15, Aug-15, Sep-15, Oct-15, Nov-15, Dec-15, Jan-16, Feb-16, Mar-16, Apr-16, May-16, Jun-16

**100 Exabytes/year by 2024!**

Science DMZ to deliver bandwidth to the end users

OSCARS for bandwidth reservation

Science DMZ

OSCARS

# ESnet: Discovery Unconstrained by Geography



SLAC at Stanford

SLAC at Stanford

*LCLS/NERSC/Esnet Superfacility demo for Photosystem II*

→ *3x network traffic*

**Network performance enables efficiency of centralized computing**

# Systems configured for data-intensive science



NERSC Cori has data partition (Phase 1, Haswell)  pre-exascale (Phase 2, KNL preproduction)
WAN-to-Cori optimized for streaming data: 100x faster from LCLS to Cori and Globus to CERN

# Real-time queue prototyped at NERSC

- In 1998 dedicated hardware; now prototype queue on Cori
- <1% of NERSC allocation
- Cryo-Em, Mass spec, Telescopes, Accelerator, Light sources



Cryo-EM: Image classification
Nogales Lab



PTF: Image subtraction pipeline



ALS: 3D Reconstruction,
rendered on SPOT web portal

# Containers for HPC Systems

- Data analysis pipelines are often large, complex software stacks
- NERSC Shifter (with Cray), supports containers for HPC systems
- Used in HEP and NP projects
    (ATLAS, ALICE, STAR, LSST, DESI)
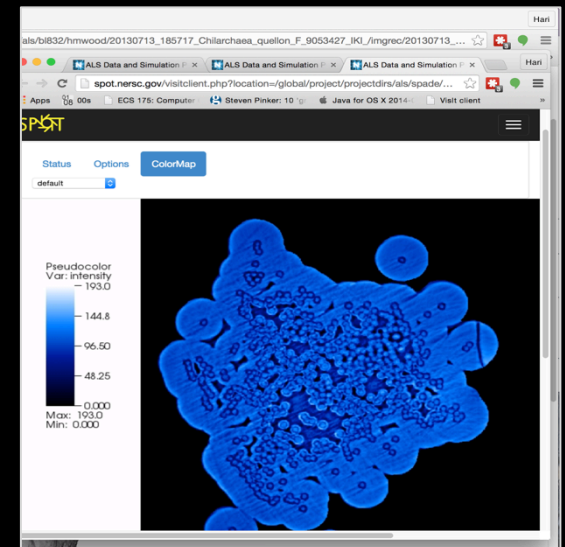


Startup Time



The Register®
Biting the hand that feeds IT

DATA CENTER   SOFTWARE   NETWORKS   SECURITY   INFRASTRUCTURE   DEVOPS   BUSINESS   HARDWA

Data Center ▸ HPC

**Cray hoists Docker containers into supercomputers**

Productivity gains without performance hits

18 Nov 2015 at 00:01, Drew Cullen

# Questions?

3. **How should undergrad/grad programs be adapted to address data challenges in future careers?**
   - New courses, (joint) majors, research institutes?
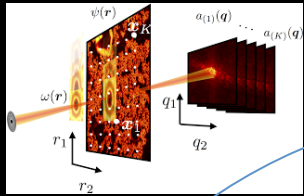
**Part 3**

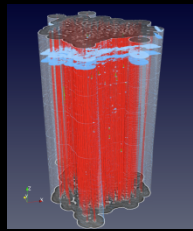# Computing, Mathematics and Statistics Research Challenges

# CAMERA: Math for the Facilities
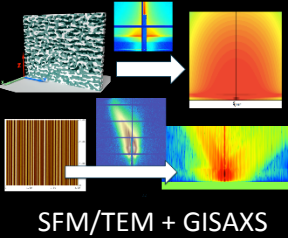## James Sethian, PI



**Designing mathematical algorithms to allow real-time analysis next to the equipment**

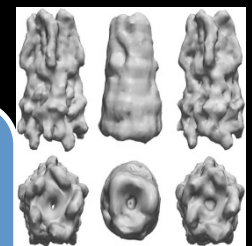Real-time streaming ptychography—ALS, delivered to NSLS2, LANL, BESSY,

**New algorithms to transform manual into automatic analysis**

Automatic image processing for the ALS/GE

SFM/TEM + GISAXS

**Multi-modal: Building the math that fuses information from multiple experiments**

CAMERA workshop on Tomography: Joint with APS, ESRF, DIAMOND, LNLS, LLNL, SSRL,....,

**Compare and integrate multiple analysis tools**

**Inventing new math and models to match new acquisition technologies**

Fluctuation scattering and single particle imaging for the LCLS

**Cultural and Sociological Challenges**

**Robust and reliable codes and data flow: workflow environments**

Workflow and access to remote supercomputers: XiCAM for ALS, SSRL, APS, NSLS2

# Analytics vs. Simulation Kernels:

| 7 Giants of Data | 7 Dwarfs of Simulation |
|---|---|
| Basic statistics | Monte Carlo methods |
| Generalized N-Body | Particle methods |
| Graph-theory | Unstructured meshes |
| Linear algebra | Dense Linear Algebra |
| Optimizations | Sparse Linear Algebra |
| Integrations | Spectral methods |
| Alignment | Structured Meshes |

# Machine Learning Mapping to Linear Algebra



Logistic Regression, Support Vector Machines

Dimensionality Reduction (e.g., NMF, CX/CUR, PCA)

Clustering (e.g., MCL, Spectral Clustering)

Graphical Model Structure Learning (e.g., CONCORD)

Deep Learning (Convolutional Neural Nets)

Sparse Matrix-Sparse Vector (SpMSpV)

Sparse Matrix-Dense Vector (SpMV)

Sparse Matrix Times Multiple Dense Vectors (SpMM)

Sparse - Sparse Matrix Product (SpGEMM)

Dense Matrix Vector (BLAS2)

Sparse - Dense Matrix Product (SpDM$^3$)

Dense Matrix Matrix (BLAS3)

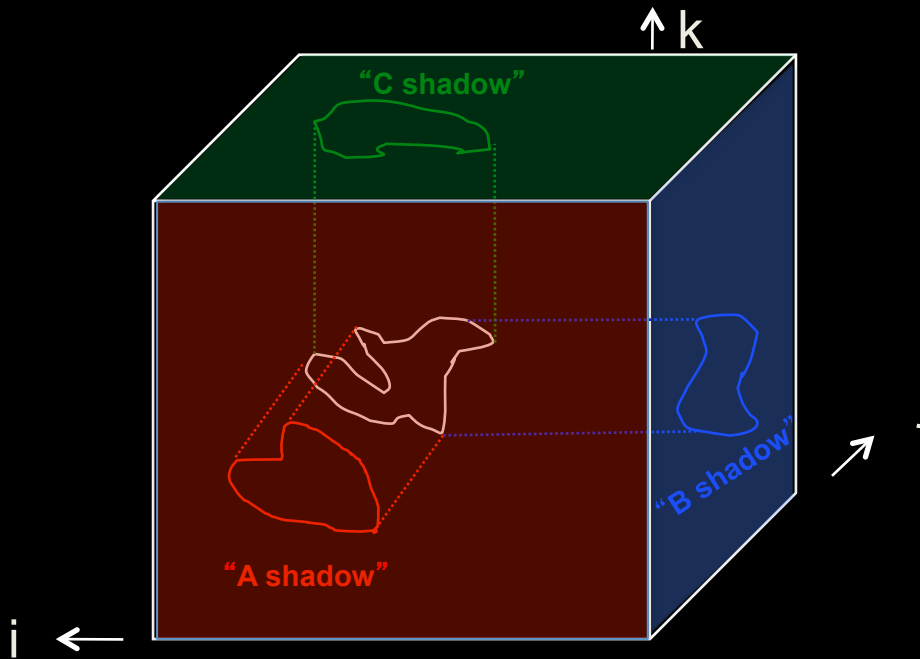Aydin Buluc, Sang Oh, John Gilbert, Kathy Yelick

# Challenge: Communication is expensive

## Communication is expensive in time and energy



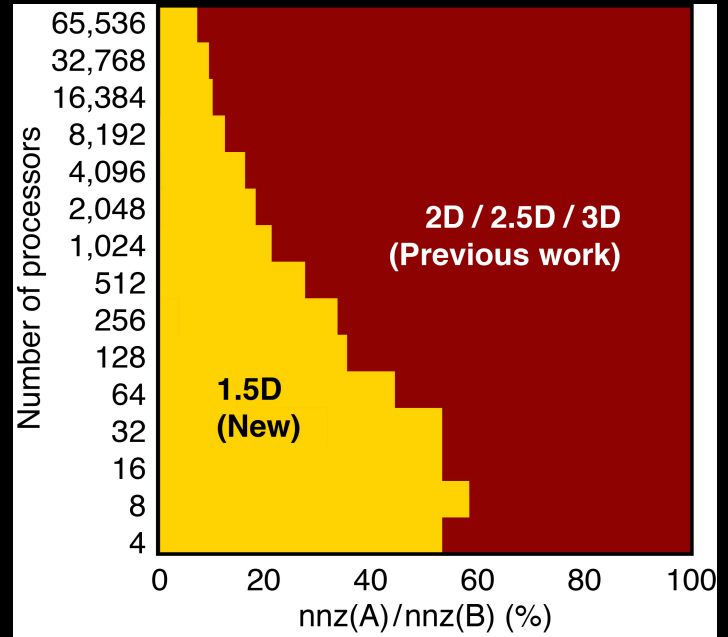## Hard to change: Latency is physics; bandwidth is money!

# Communication-Avoiding Algorithms



Matrix Multiplication code has a 3D iteration space; each point is a */+

```
for i
  for j
    for k  C[i,j] …  A[i,k] …  B[k,j]  …
```

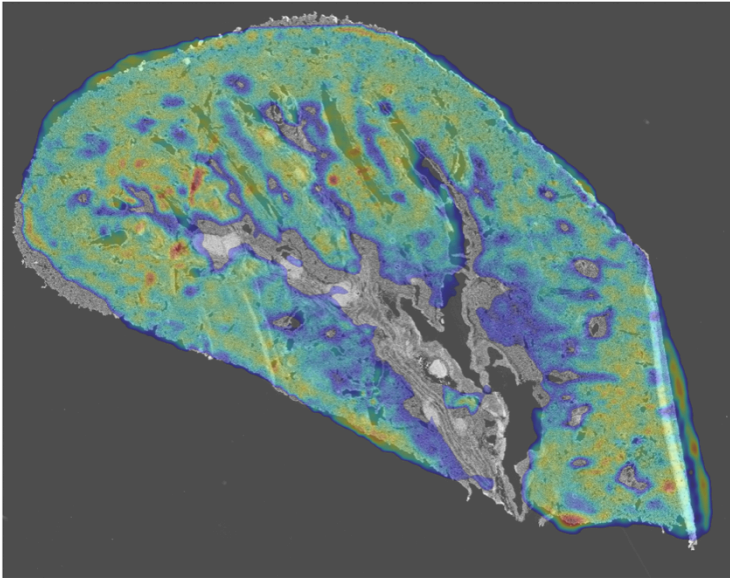Model for choosing communication-optimal algorithms for sparse matrices

Demmel et al on LA; Christ et al generalization

Koanantakool & Yelick

# Interactive Analytics using Jupyter



**Science notebooks through Jupyter (iPython)**

- Widely used in science
- Interactive HPC LDRD

**Deployed at NERSC:**

- >100 users pre-production

*Fernando Perez et al*

# Random Access Analytics

- **Genome assembly "needs shared memory"**

**Global Address Space**



**Scales to 15K+ cores**

**4 minutes for human**

**First ever solution**

Distributed hash table
- Low overhead communication
- Remote atomics, caching
- Locality-aware hashing

*E. Georganas, A. Buluc, J. Chapman, S. Hofmeyr, C. Aluru, R. Egan, L. Oliker, D. Rokhsar, K. Yelick*
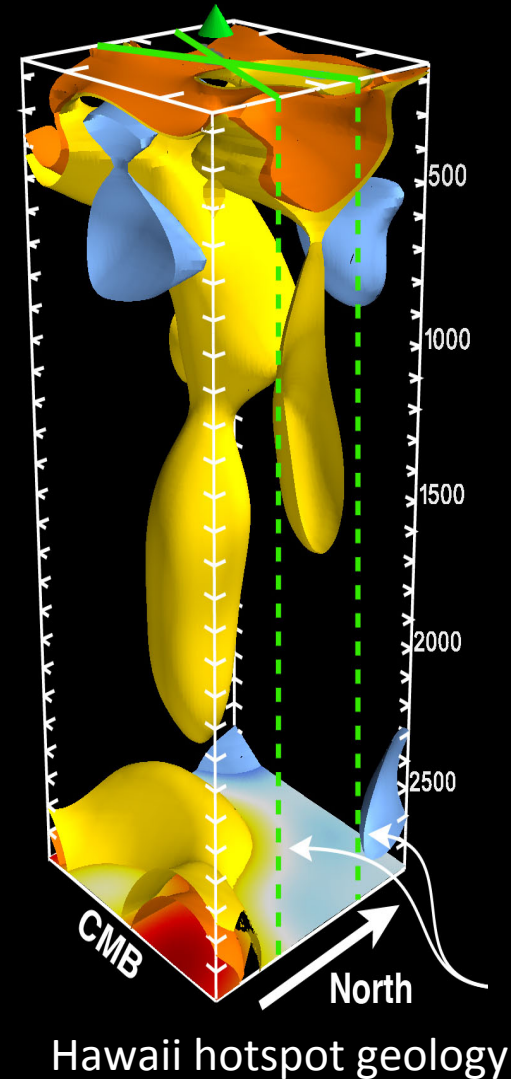
# Data Fusion for Observation with Simulation



- **Unaligned data from observation**
- **One-sided strided updates**

Scott French, Y. Zheng, B. Romanowicz, K. Yelick
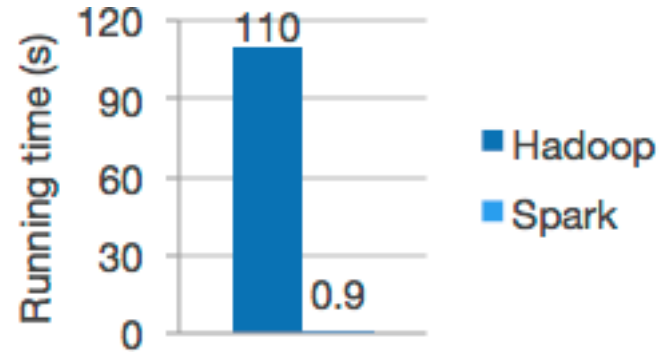


Hawaii hotspot geology
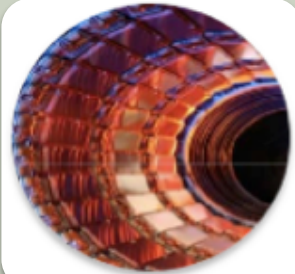
# Productive Programming



**Speed**
Run programs up to 100x faster than Hadoop
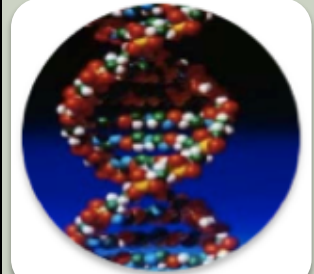MapReduce in memory, or 10x faster on disk.

- **High failure rate**
- **Slow network**
- **Fast (local) disk**

**And Spark is still 10x+ slower than MPI**

# Architectures for Data vs. Simulation



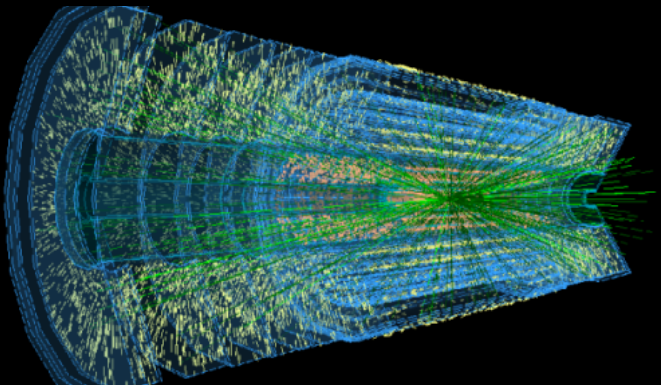**Massive Indepen-dent Jobs for Analysis and Simulation**

**Random access, large data Analysis**

**Different architectures for simulation?  Can simulation use data architectures?**
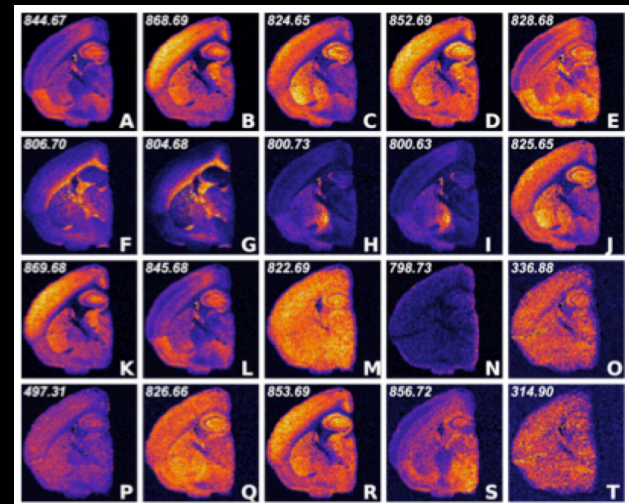
# Data processing with special purpose hardware

- General trend towards specialization for continued performance growth
- Data processing (on raw data) will be first in DOE



Particle Tracking with Neuromorphic chips

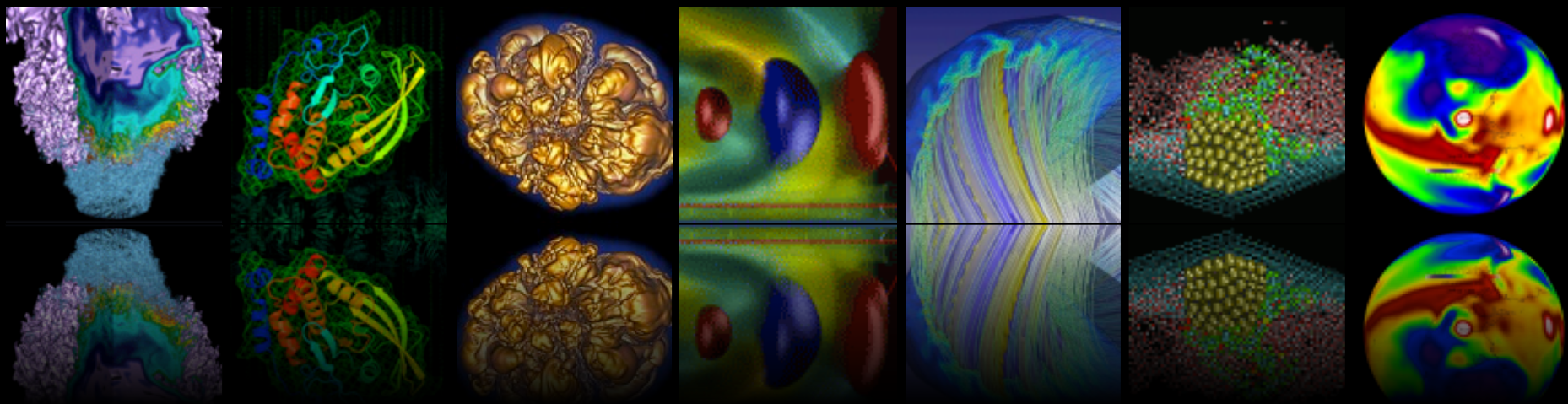Computing in Detectors



Deep learning processors for image analysis

FPGAS for genome analysis

# Questions?

4. Are there open problems or expertise gaps in computing/math/stat/data be addressed?

# Questions?

1. Are there MSU examples of "science at the boundary" of simulation and observation?
   – How should you take advantage of these opportunities?
2. What are the largest and most complex sources of research data at MSU?
   – What types of data/CS/math/stat challenges are there?
3. How should undergrad/grad programs be adapted to address data challenges in future careers?
   – New courses, (joint) majors, research institutes?
4. Are there open problems or expertise gaps in computing/math/stat/data be addressed?

# Extreme Data Science

The scientific process is poised to undergo a radical transformation based on the ability to access, analyze, simulate and combine large and complex data sets.

Slides: http://www.cs.berkeley.edu/~yelick/talks