

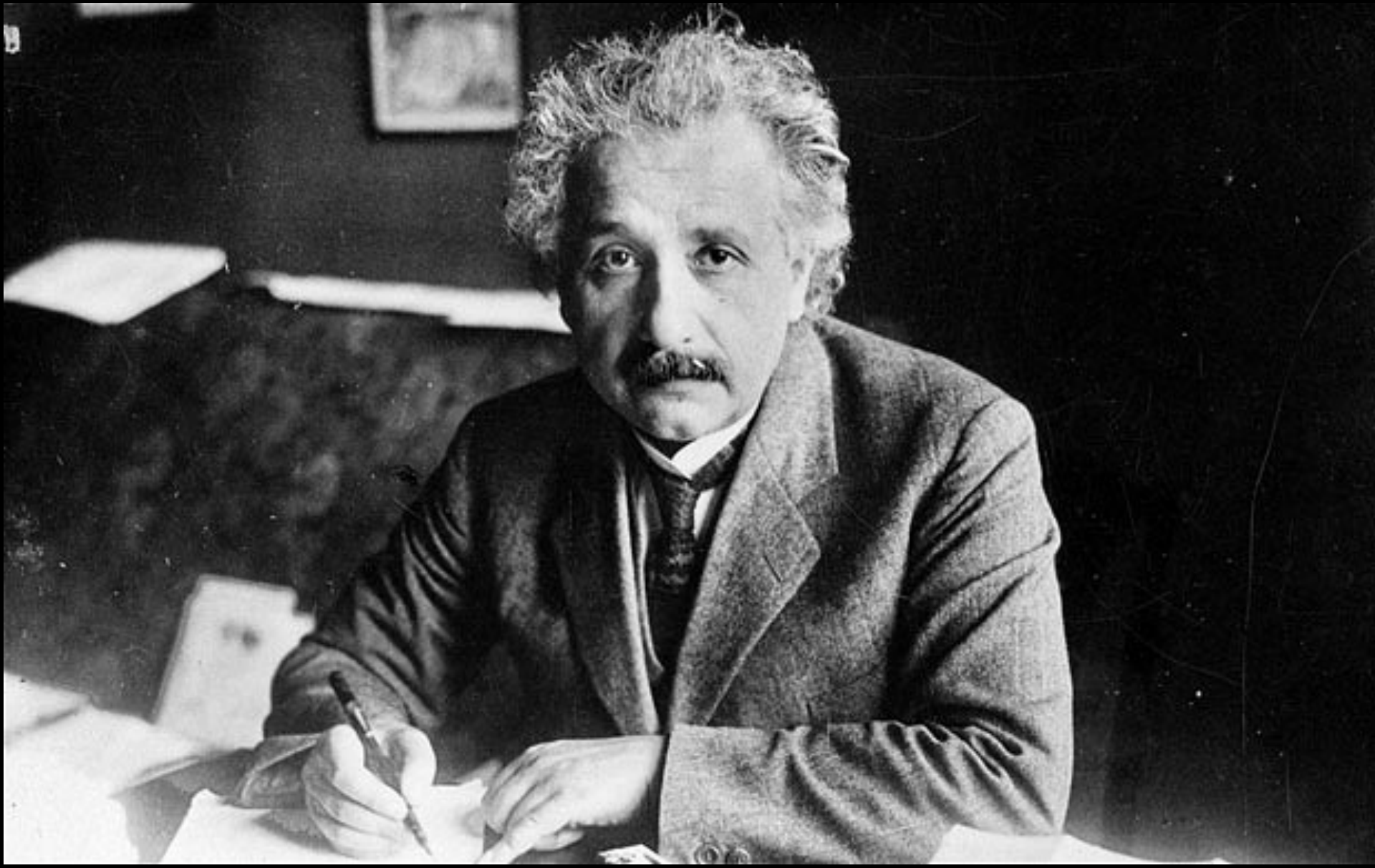
More Data, More Science... and Moore's Law?

Kathy Yelick

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

Science is poised for transformation

Old School Scientists: The Lone Scientist



The Legacy of Team Science



Radiation Lab staff on the magnet yoke for the 60-in cyclotron, 1939, including:

E. O. Lawrence

Edwin McMillan

Luis Alvarez

J. Robert Oppenheimer

Robert R. Wilson

New Scientists



17-year-old Brittany Wegner creates breast cancer detection tool that is 99% accurate on a minimally invasive, previously inaccurate test.

Machine Learning + Online Data + Cloud Computing

Experimental Science is Changing

By using our website you agree to our use of cookies in accordance with our cookie policy.

OK

PRIVACY POLICY +

JAX[®] MICE & SERVICES



JAX[®] Mice are the highest quality and most-published mouse models in the world. Take advantage of our large inventories of common inbred strains and the convenience of having your breeding and drug efficacy needs met by the leading experts in mouse modeling.

Search for Mice

Advanced Mice Search

Search for mice by strain, stock, gene, allele and synonyms



Breed Your Mouse

Test Your Drug

Cryopreserve Your Mouse

Computing Sciences at Berkeley Lab



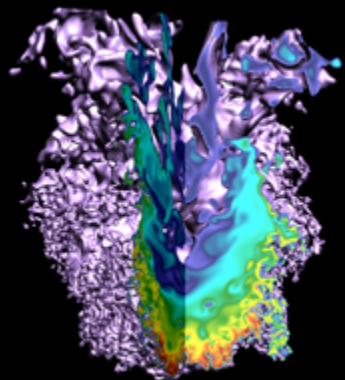
NERSC: State-of-the art supercomputing for the broad science community – over 7000 users, 700 applications mostly in simulation

High Performance Computing in Science

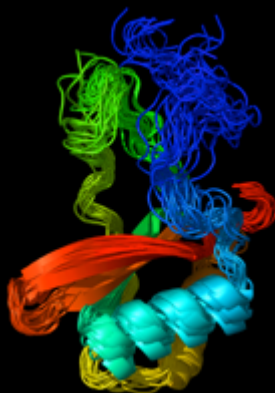
Computers are used to understand things that are



Too Big



Too fast



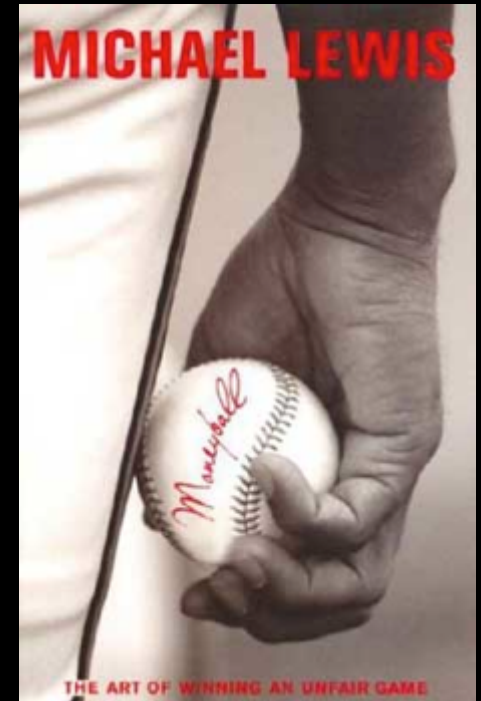
Too Small



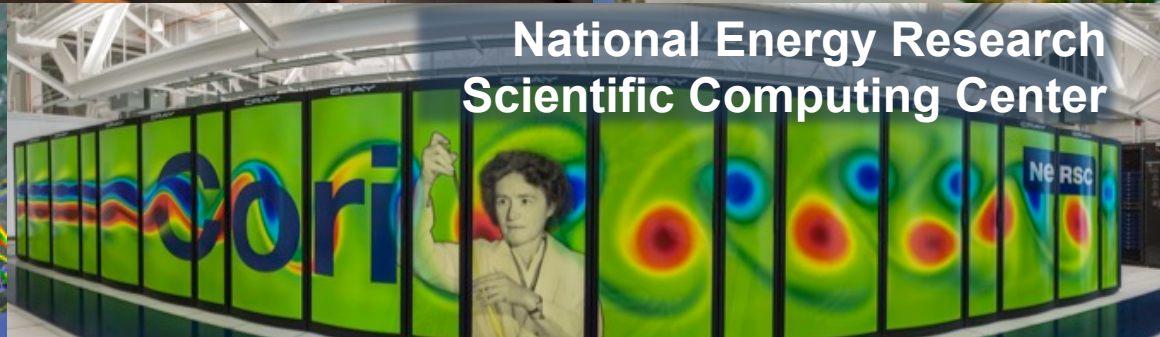
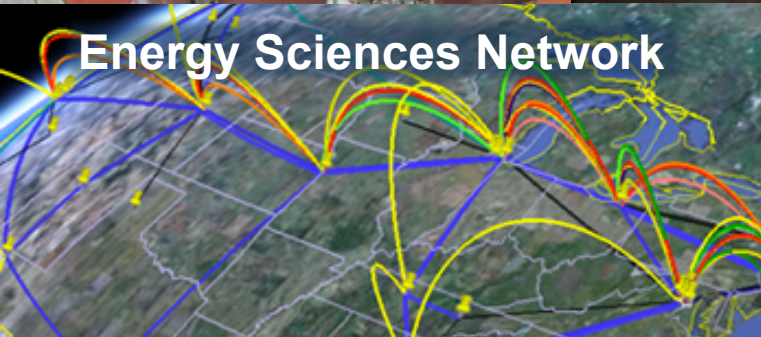
Too slow

for experiments alone, so simulations are used

“Big Data” Changes Everything...What about Science?



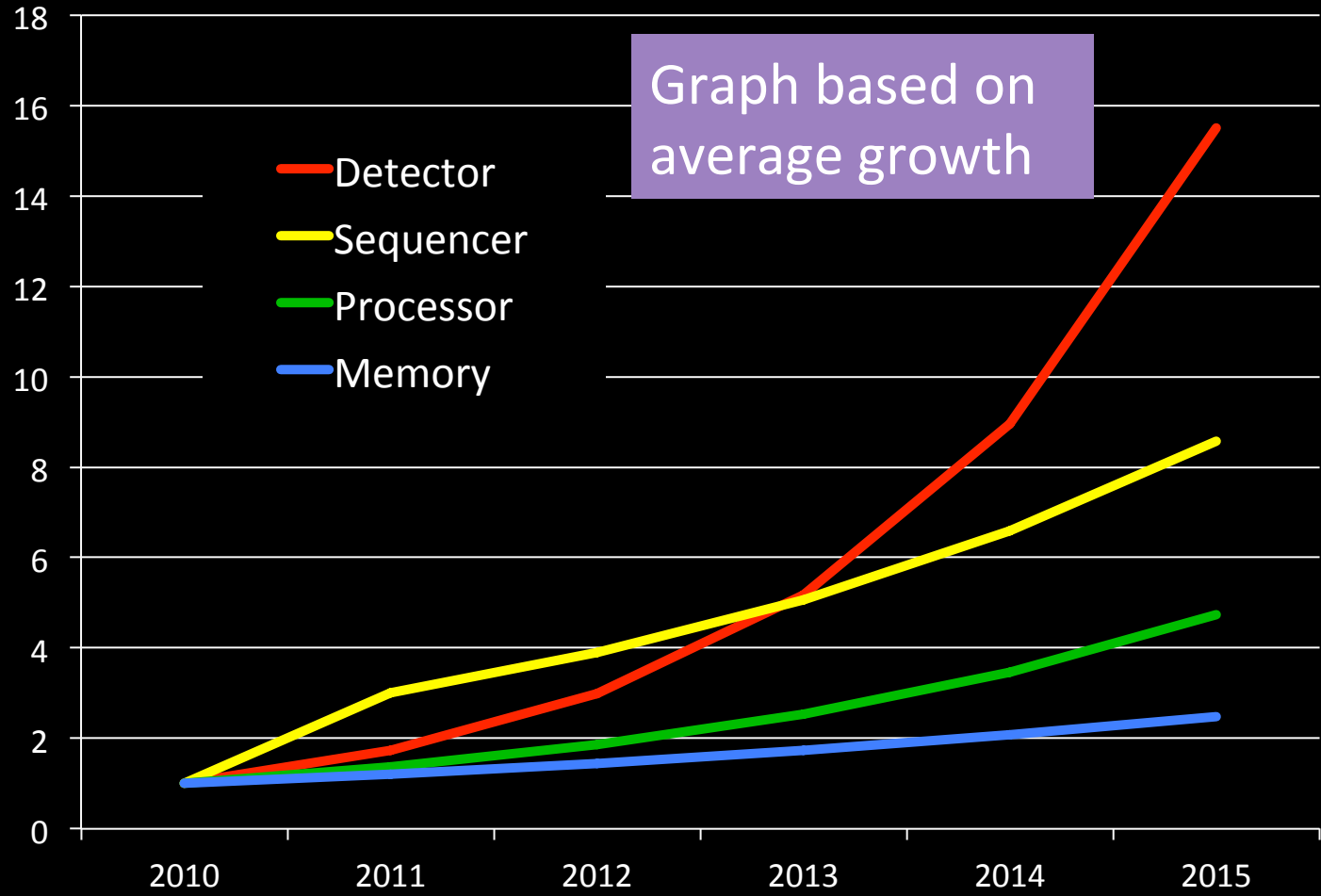
Berkeley Lab's Advanced Facilities Enable World-Leading Science



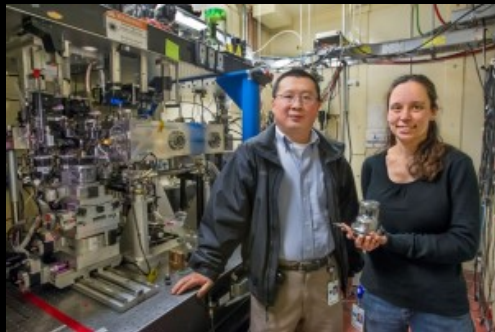
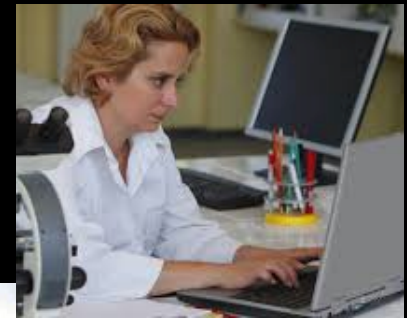
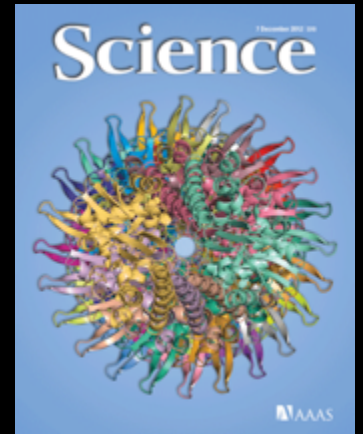
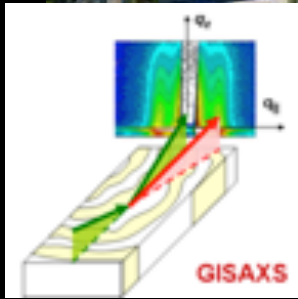
About 10,000 visiting scientists (~2/3 from universities) use Berkeley Lab research facilities each year, which provide some of the world's most advanced capabilities in materials science, biological research, computation and networking

Data Growth is Outpacing Computing Growth

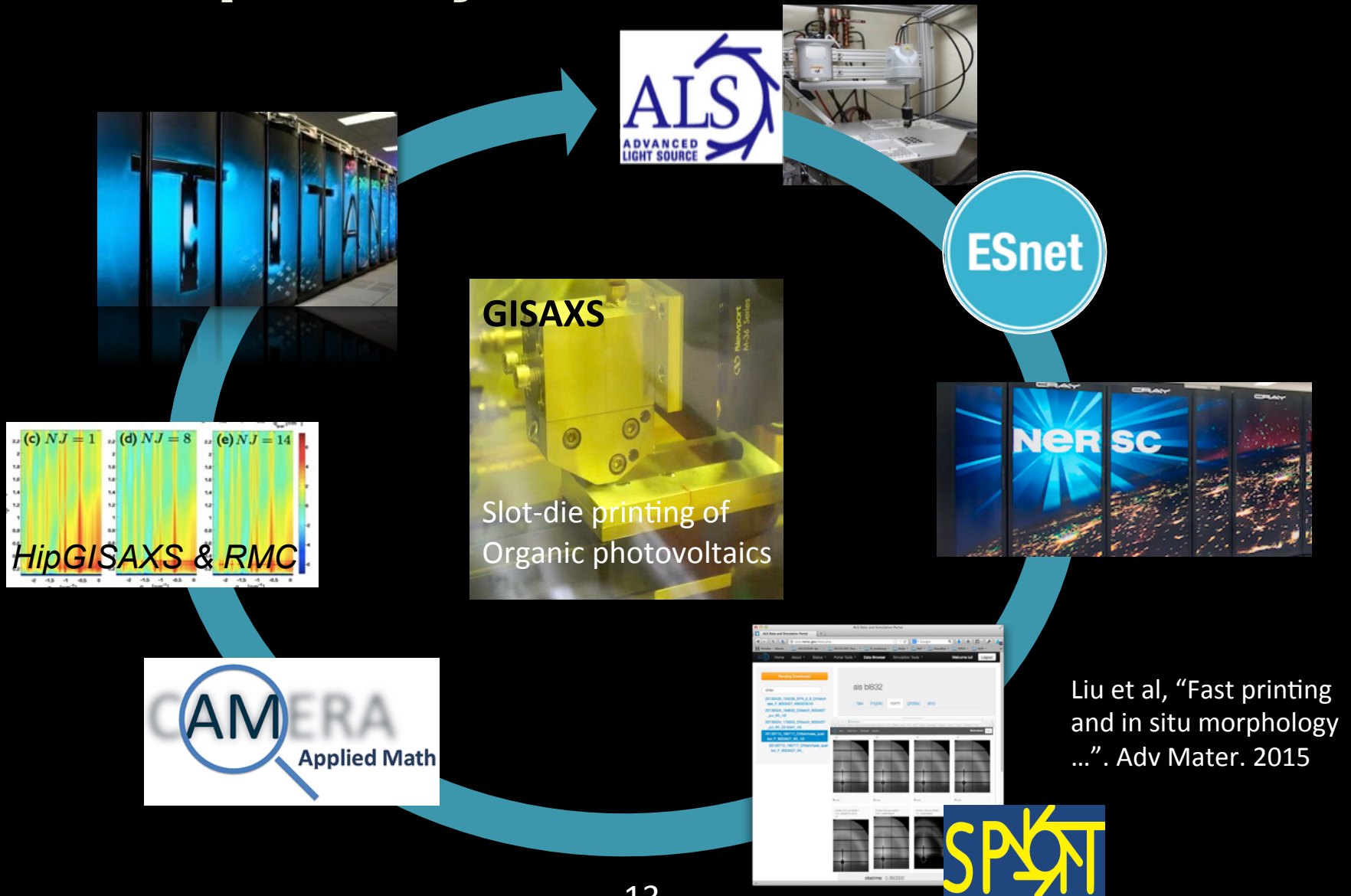
Projected Data Rates Relative to 2010



Old School Scientific Workflow

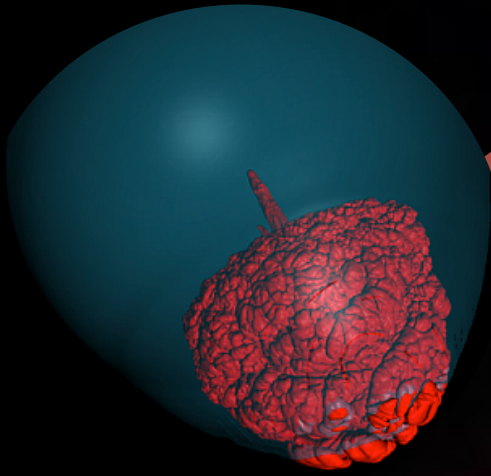


Computing, experiments, networking and expertise in a “Superfacility” for Science



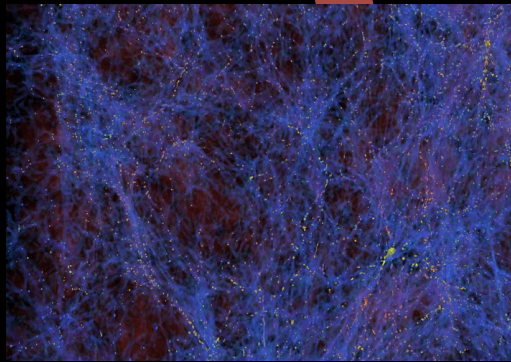
**Science at the boundary of theory
and experiment ... simulation and
data analytics**

Integration of Simulation and Observational Science



Intermediate Palomar
Transient Factory with DESI,
CMB-S4 and LSST coming

A. Goobar, P. Nugent, et al
(2017) Science



Simulations aid in
interpreting data

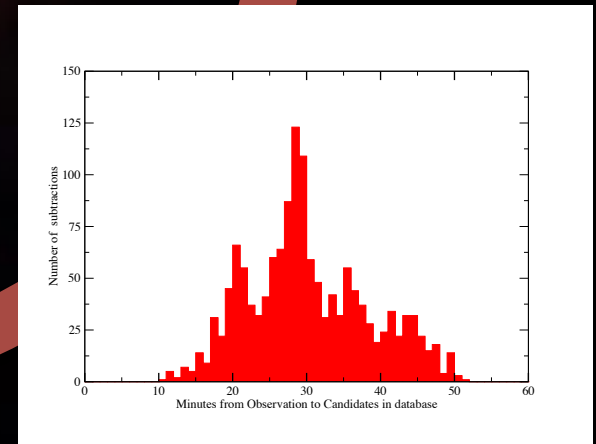
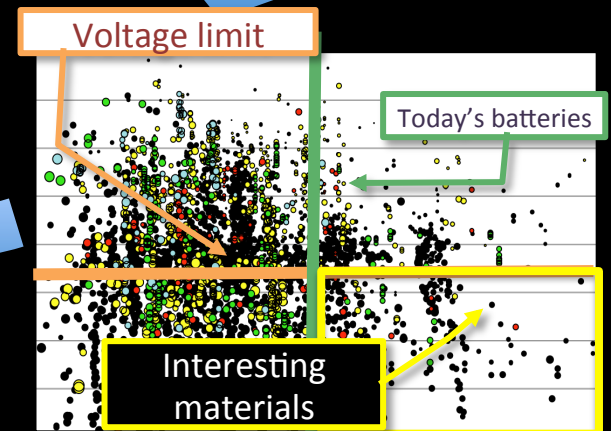


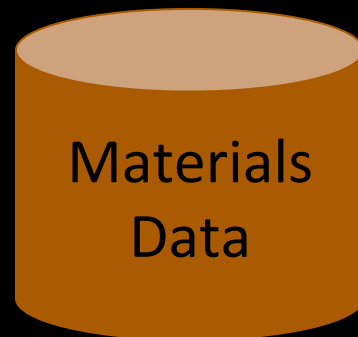
Image subtraction, machine
learning in minutes

Re-Use and Re-Analyze Previously Collected Data

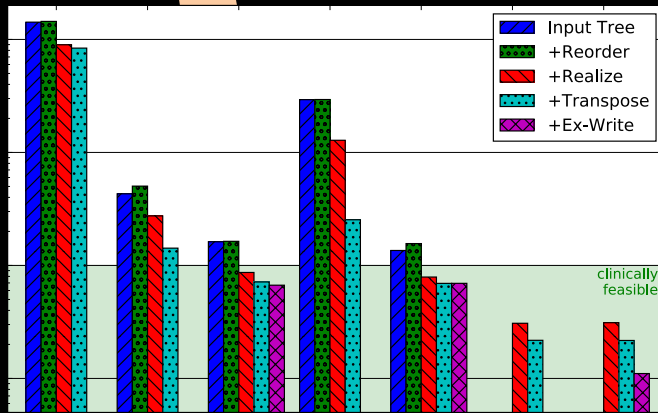
- **Materials Genome Initiative**
 - Materials Project: Over 10,000 users!
 - “World-Changing Idea of 2013”



Computers programs run by “bots”

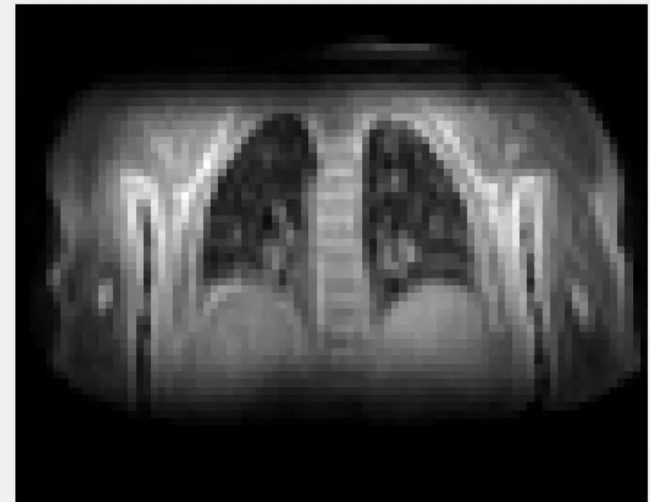


Real-Time Analytics in Health



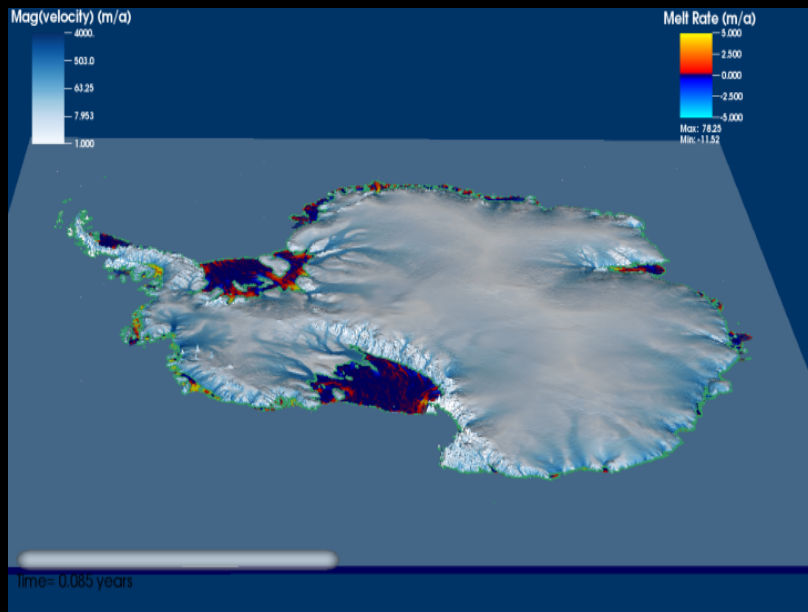
3 min goal (1 sec/iteration)

Michael Driscoll HPC optimization



Compressed Sensing Approach by Mike Lustig et al
MRI results Wenwen Jiang

Data and Simulation in the environment



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

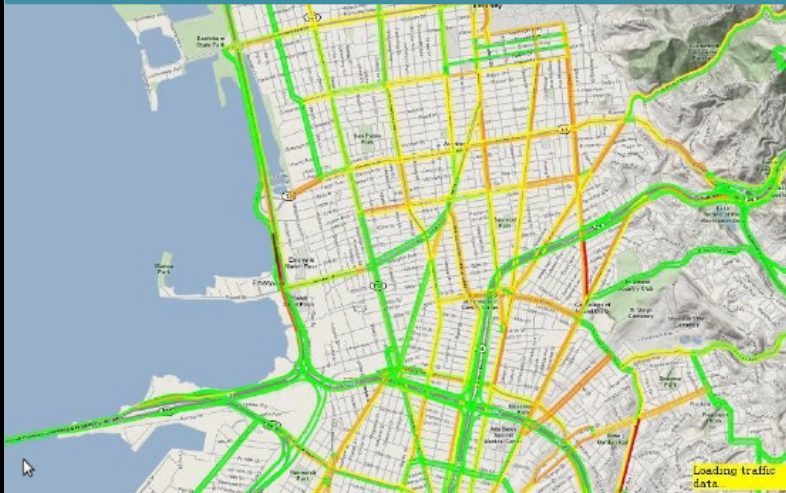


Genomes to watersheds Scientific Focus Area

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.

Science in embedded sensors

Transportation Modeling



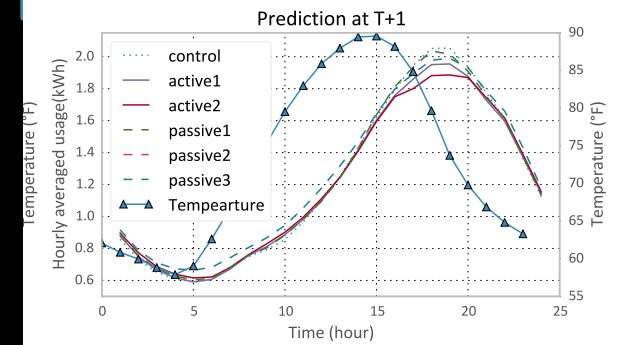
Power Grid Modeling



Scenario Prediction, Planning

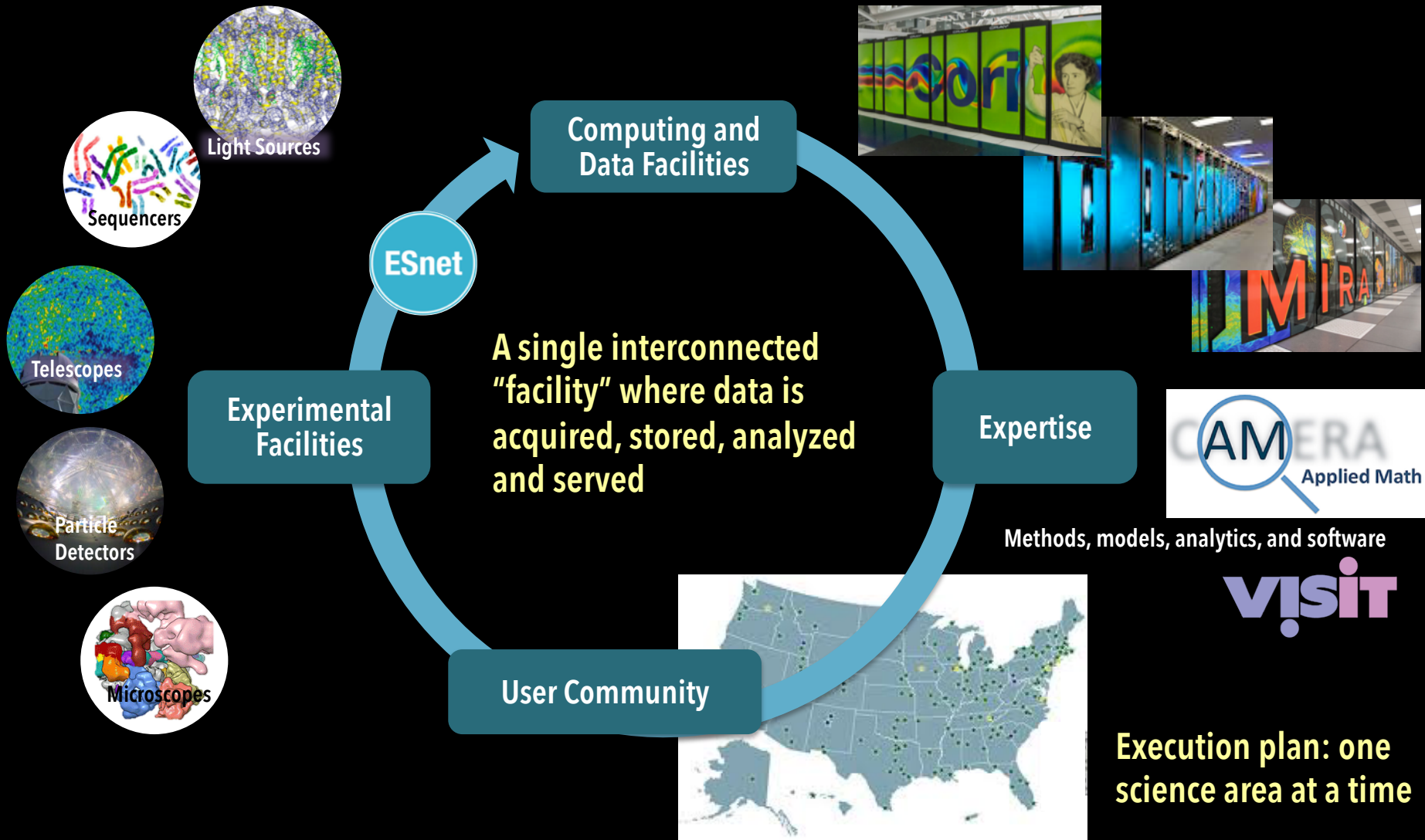


Decision Science

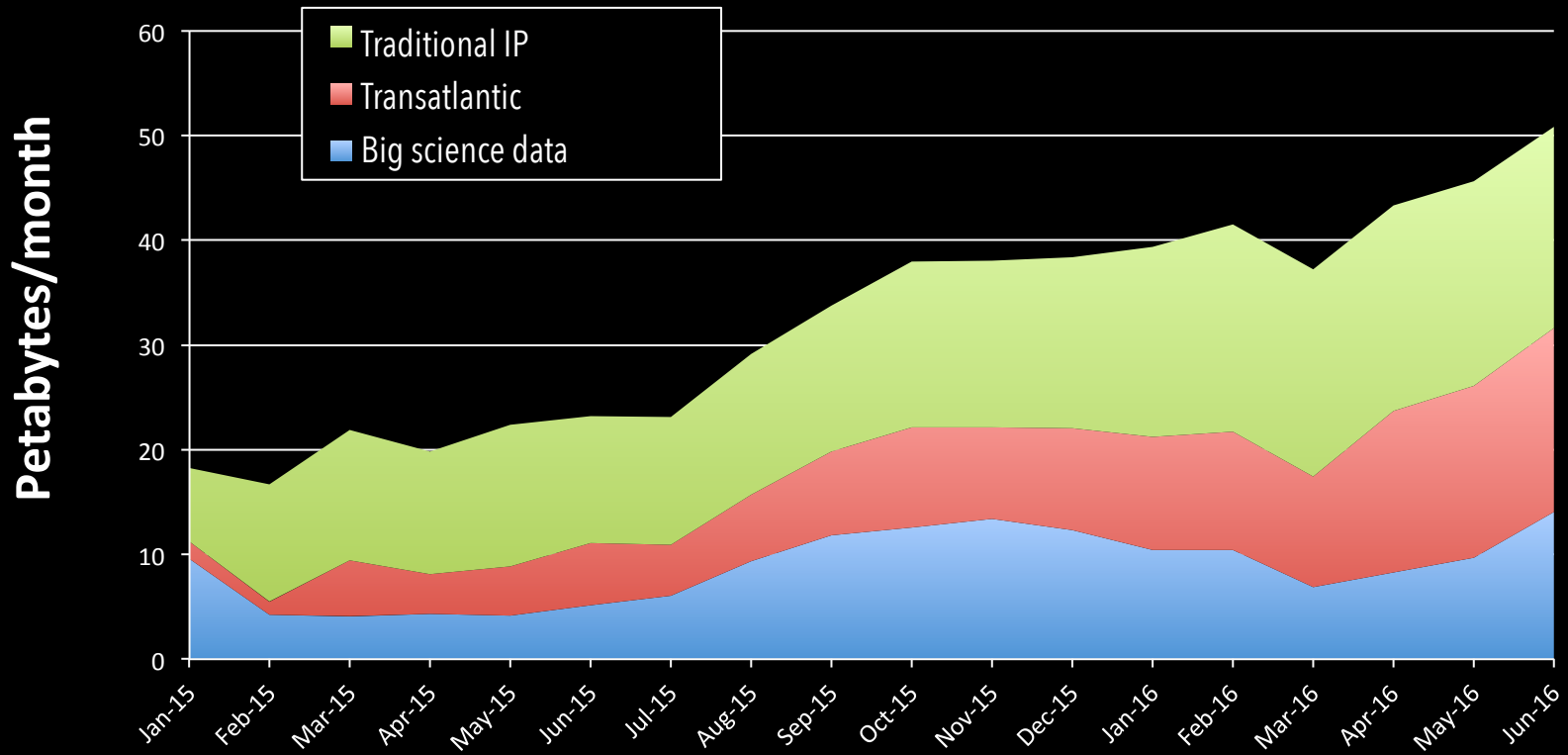


**Computing and network facilities
need to adapt**

Superfacility: Integrated network of experimental and computational facilities and expertise



ESnet: Data driven science drives network capacity



Science DMZ to deliver bandwidth to the end users

OSCARS for bandwidth reservation

100 Exabytes/year by 2024!



Systems configured for data-intensive science



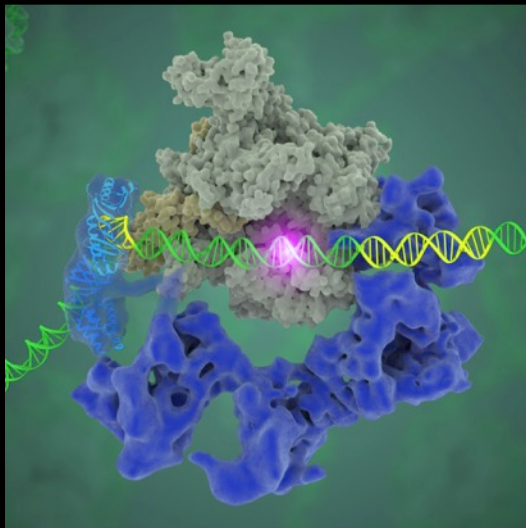
NERSC Cori has data partition (Haswell) and pre-exascale (KNL)

NVRAM file system with close to 2 PB at 2 TB/sec

WAN-to-Cori optimized for streaming data: 100x faster from LCLS to Cori and Globus to CERN

Real-time queue prototyped at NERSC

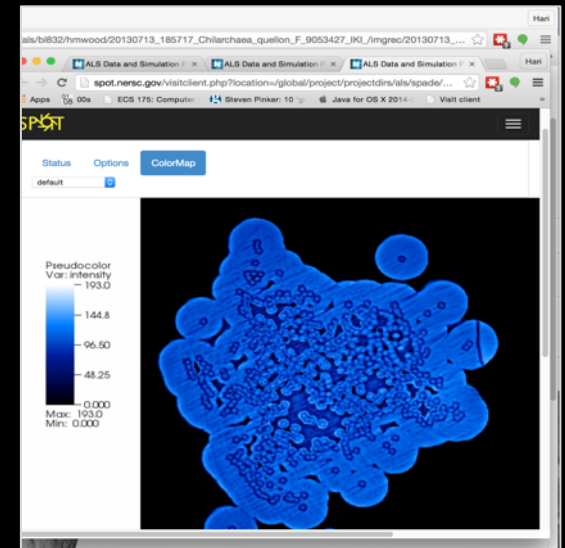
- In 1998 dedicated hardware; now prototype queue on Cori
- <1% of NERSC allocation
- Cryo-Em, Mass spec, Telescopes, Accelerator, Light sources



Cryo-EM: Image classification
Nogales Lab



PTF: Image subtraction pipeline



ALS: 3D Reconstruction,
rendered on SPOT web portal

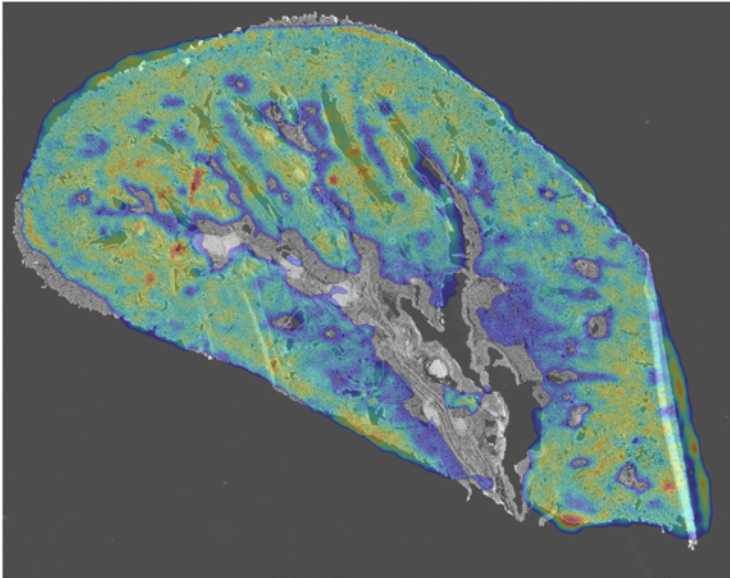
Interactive Analytics using Jupyter

```
In [10]: # overlaying the small H&E and MS images

registered_ms_image = ird.transform_img_dict(my_images[2], result)
big_registered_ms_image = imresize(registered_ms_image, optical_image.shape, interp='bicubic')

# cut out low intensity region of MS image for easy viewing of underlying H&E
masked_big_ms_image = np.ma.masked_where(big_registered_ms_image < 100, big_registered_ms_image)

# plot the two images overlaid
f = plt.figure(1, figsize=(20, 20))
plt.imshow(optical_monochrome, alpha=0.7, cmap=cm.Greys_r)
plt.imshow(masked_big_ms_image, alpha=0.3, cmap=cm.jet)
plt.axes().set_axis_off()
```



Science notebooks through Jupyter (iPython)

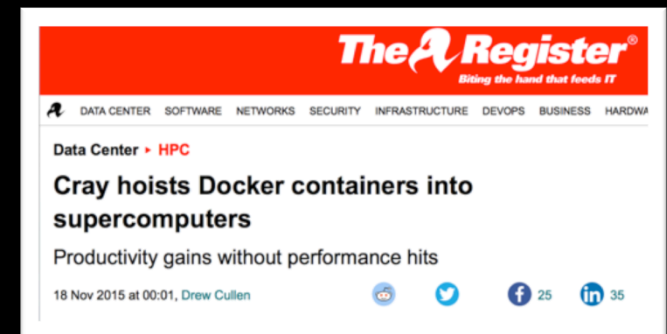
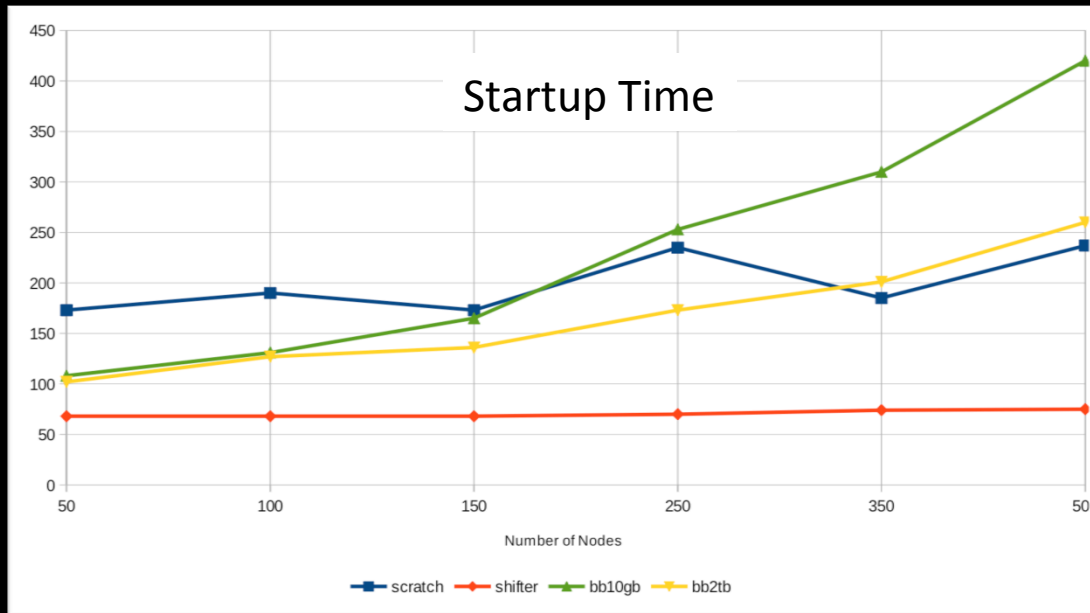
- Widely used in science
- Interactive HPC LDRD

Deployed at NERSC:

- >100 users pre-production

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)



Old School Scientific Data Search

Safari File Edit View History Bookmarks Window Help

www.google.com/search?tbs=sbi:AMhZZIu-Ft1o4xXIjhVjclUv_1GtY_1M9gV_1hy

Berkeley Lab (...) TeamSnap :: M... Google CalMail - You... Search Results...

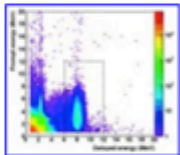
+You Search **Images** Mail Drive Calendar Sites Groups More -

CalMail - You must be logged in to a page.

Google Antineutrinos.jpg

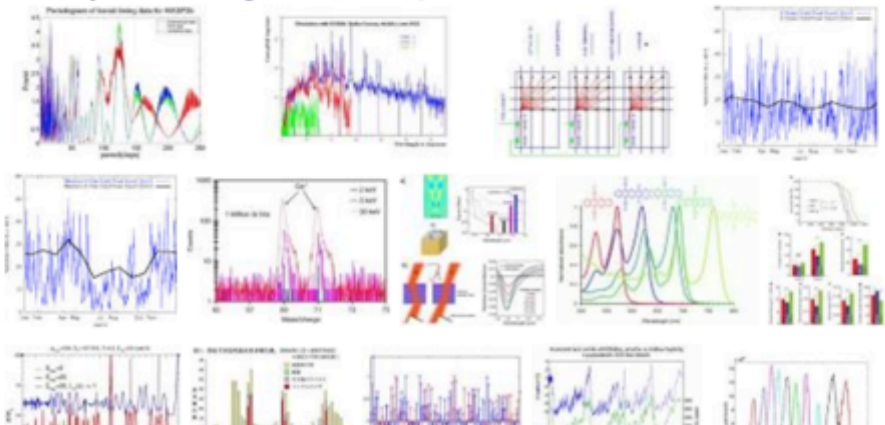
Web **Images** Maps Shopping More Search tools

Tip: Try entering a descriptive word in the search box.

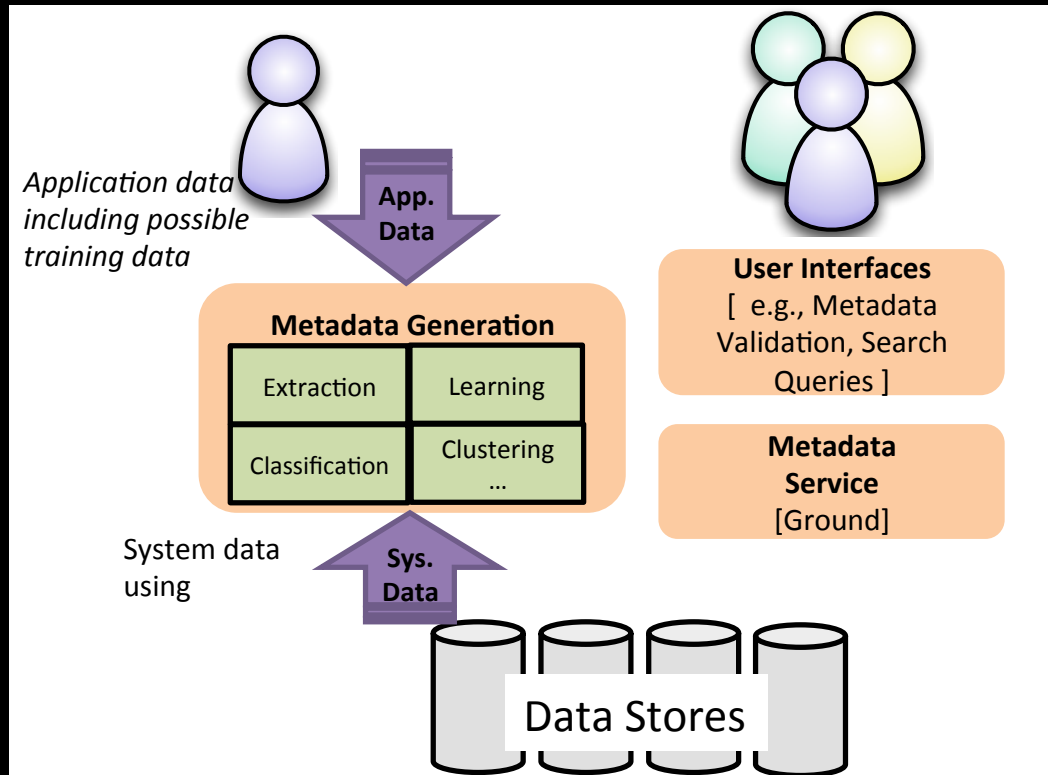
 Image size:
153 × 133

No other sizes of this image found.

[Visually similar images](#) - Report images



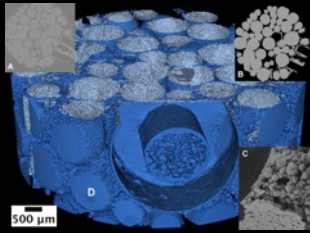
Automated Search, Meta-Data Analysis, and On-Demand Simulation



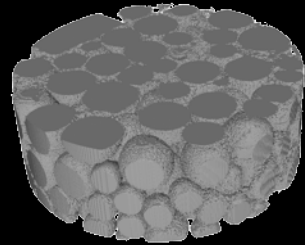
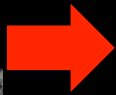
Automated metadata extraction using machine learning

**Computational research
challenges are substantial**

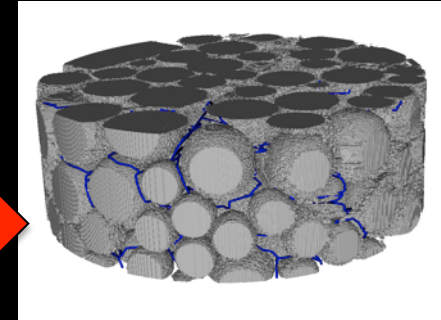
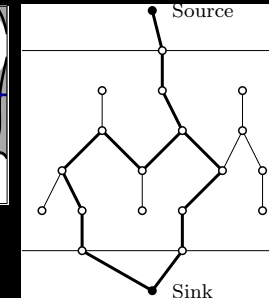
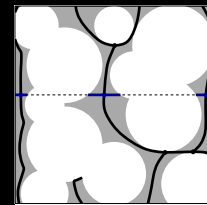
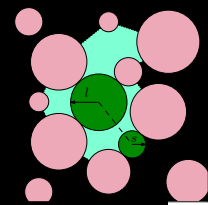
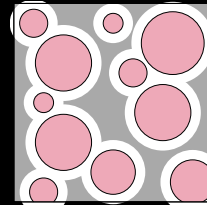
Software implementations at scale in pipeline



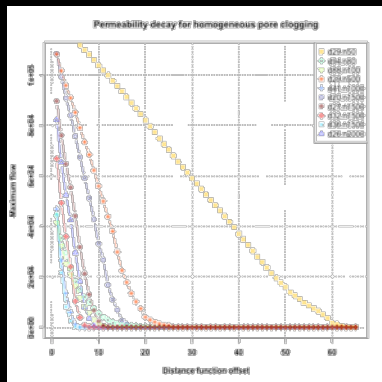
MicroCT
imaging



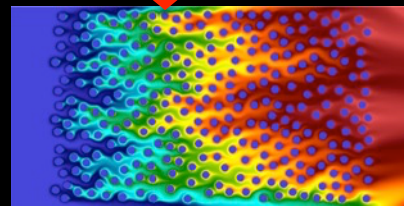
Segmentation



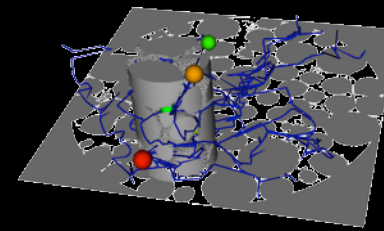
Topological
Analysis



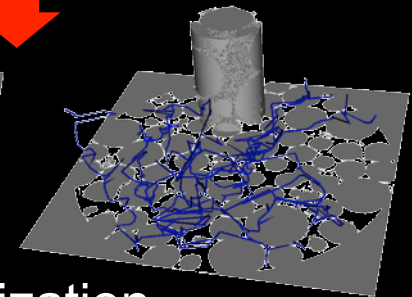
Analysis



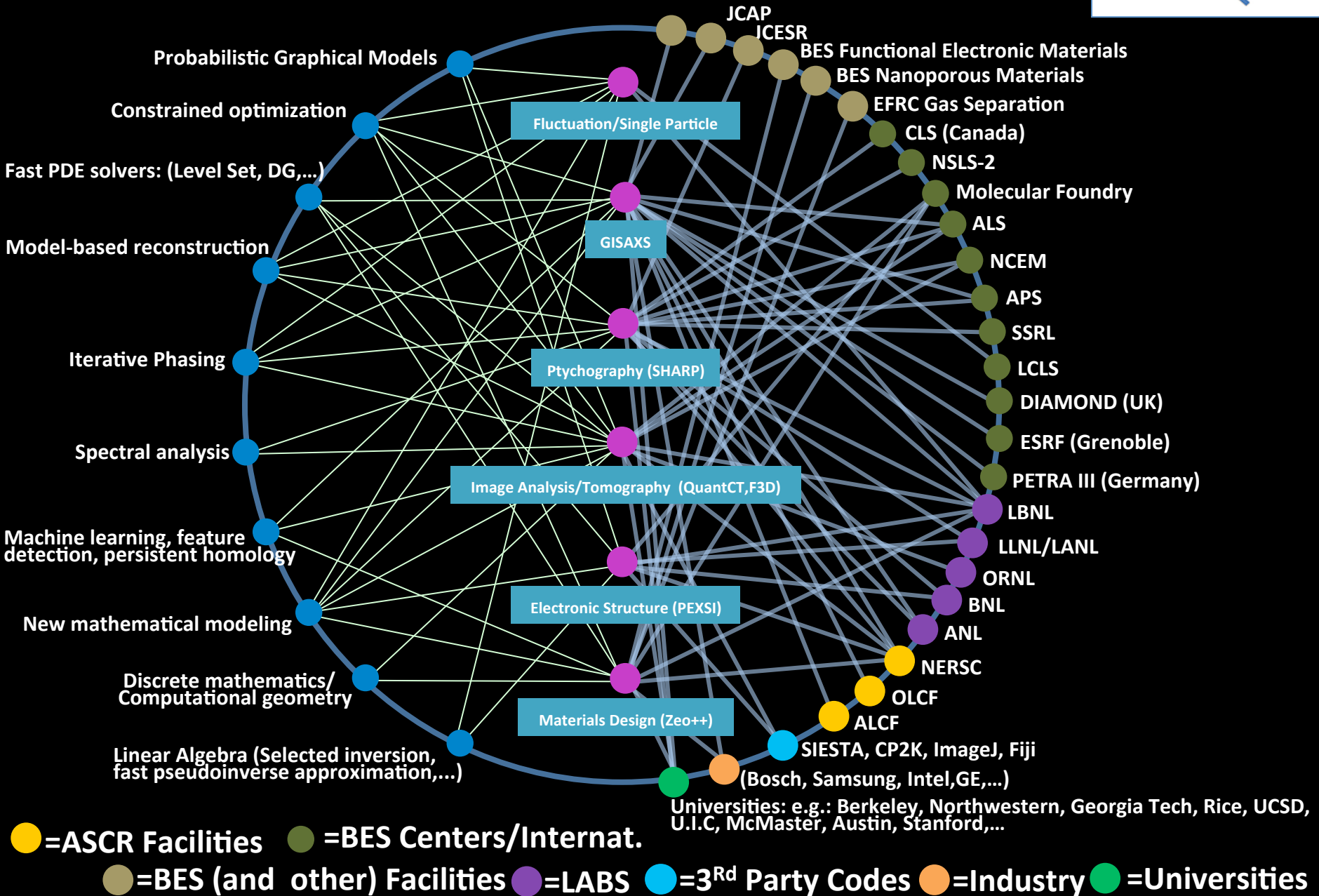
Simulation



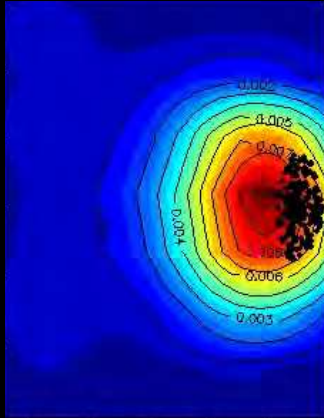
Visualization



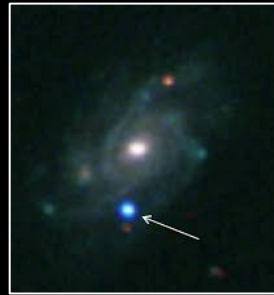
CAMERA: Mathematics for Facilities



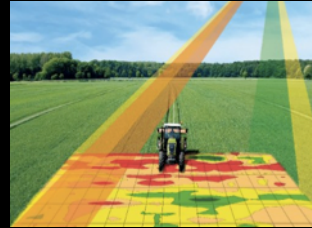
Machine Learning for Science



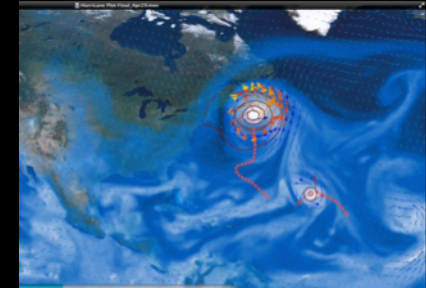
Accelerators



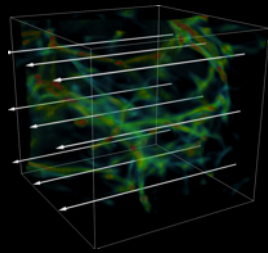
Images in cosmology,
light sources, etc.



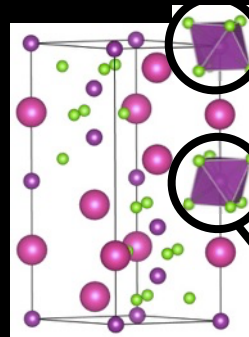
Biology



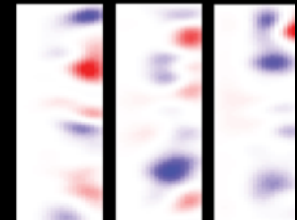
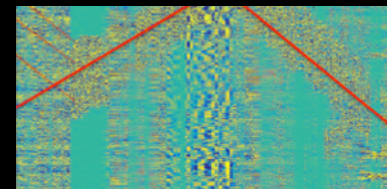
Climate



Cosmology simulation



Chemistry



Data Complexity

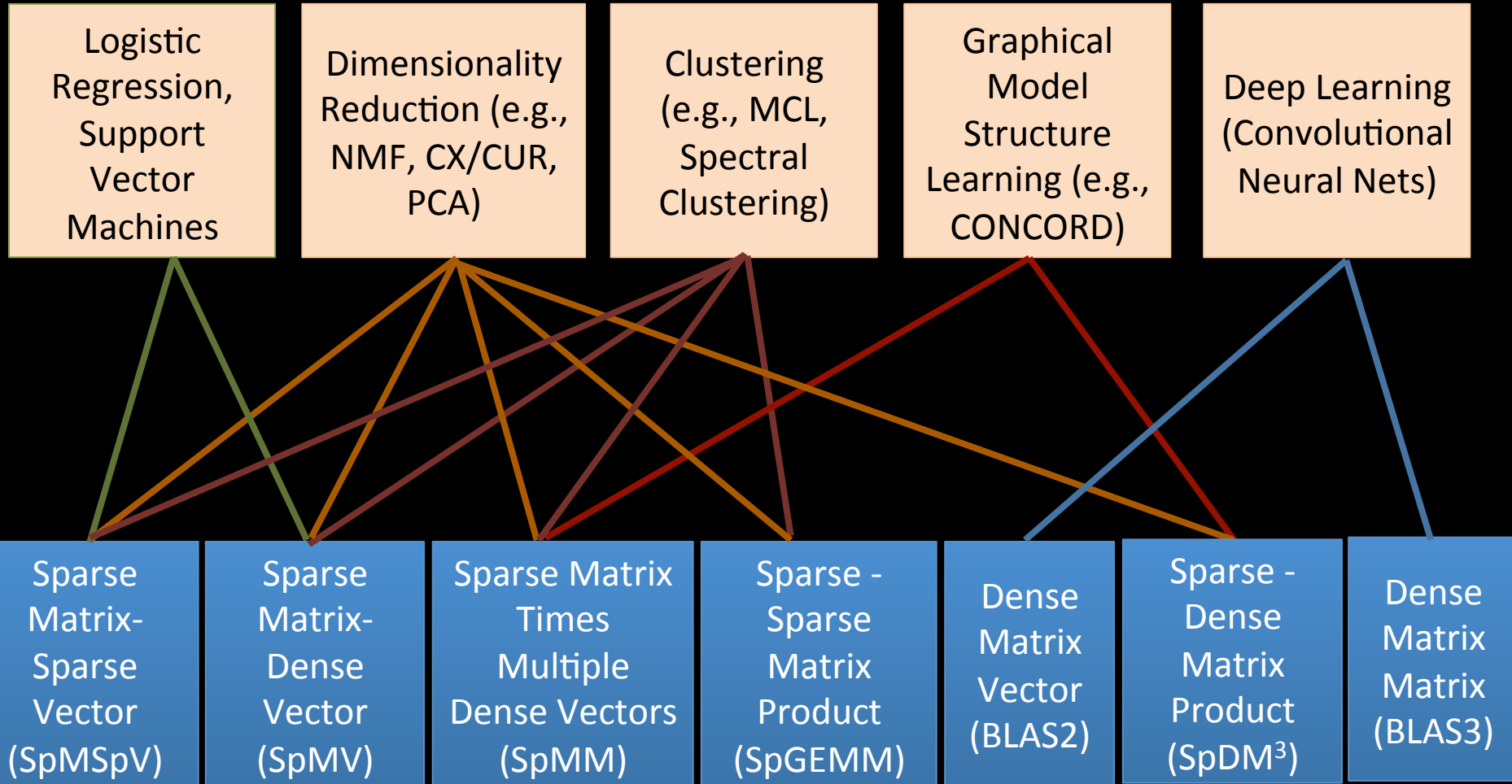
Interpretability

Performance and Scale

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 Sparse Linear Algebra
Optimizations	
Integrations	Spectral methods
Alignment	Structured Meshes

Machine Learning Mapping to Linear Algebra



Random Access Analytics

- Genome assembly “needs shared memory”

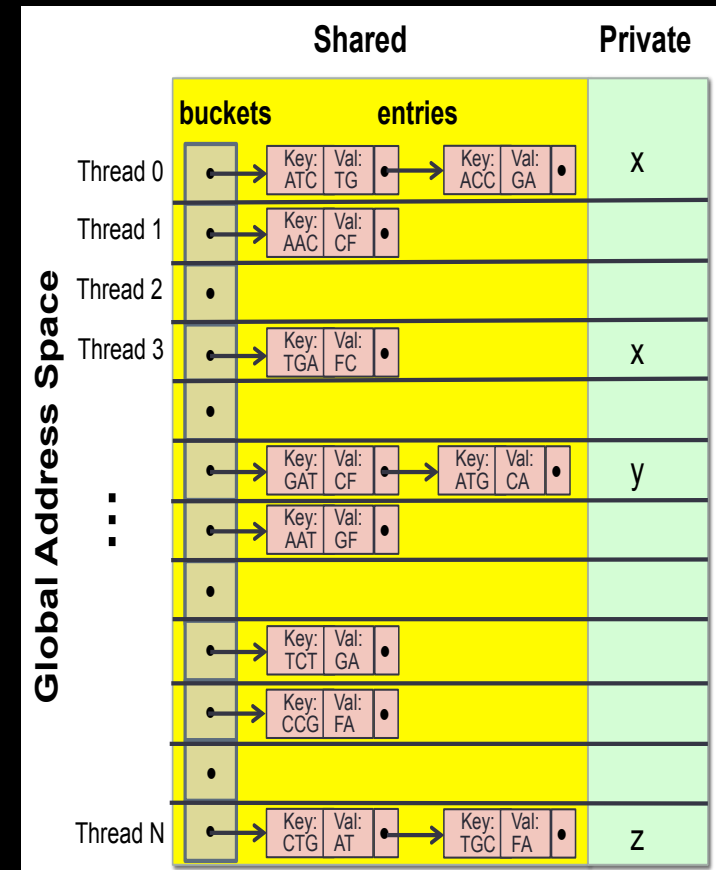
Global Address Space

- Low overhead communication
- Remote atomics
- Partitions for any structure

Scales to 15K+ cores

Under 10 minutes for human

First ever solution

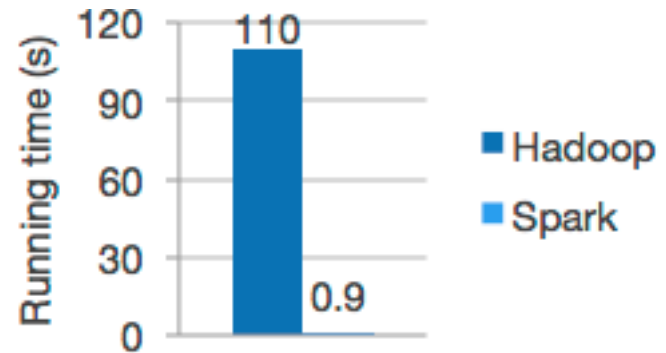


Productive Programming



Speed

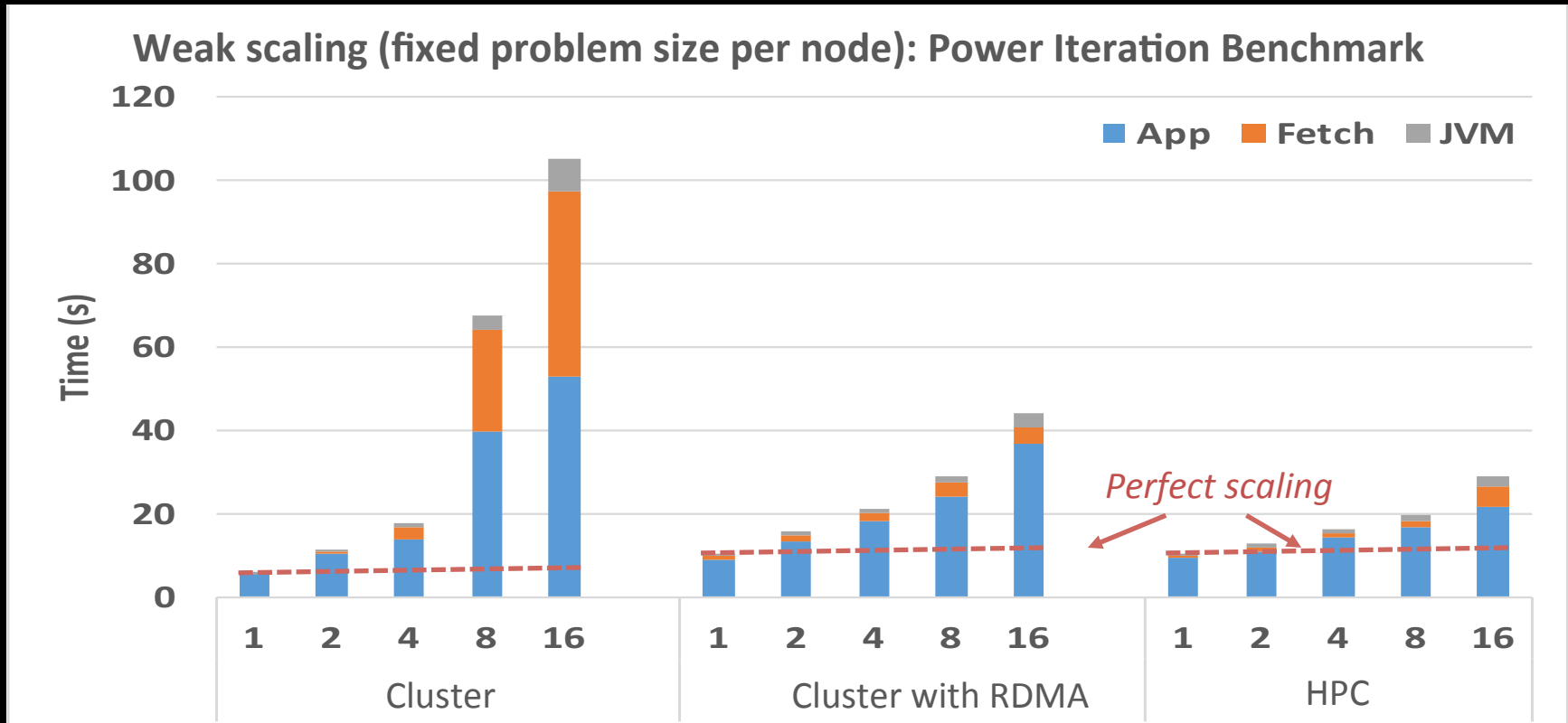
Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.



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

And Spark is still 10x+ slower than MPI

SPARK Analytics on HPC



SPARK on HPC vs. clusters

- Network, I/O, and virtualization all key to performance
- Increased scale from O(100) to O(10,000) cores

Filtering, De-Noise and Compressing Data

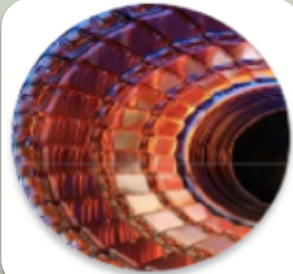


AmeriFlux & FLUXNET: 750 users access carbon sensor data from 960 carbon flux data years

Arno Penzias and Robert Wilson discover Cosmic Microwave Background in 1965

**How will we get enough
computing for these problems?**

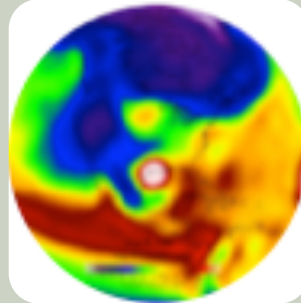
Architectures for Data vs. Simulation



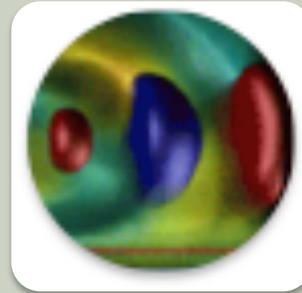
**Separate
Jobs**



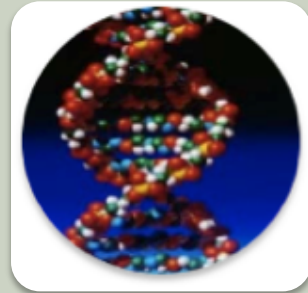
**Compute
Intensive**



**Nearest
Neighbor**



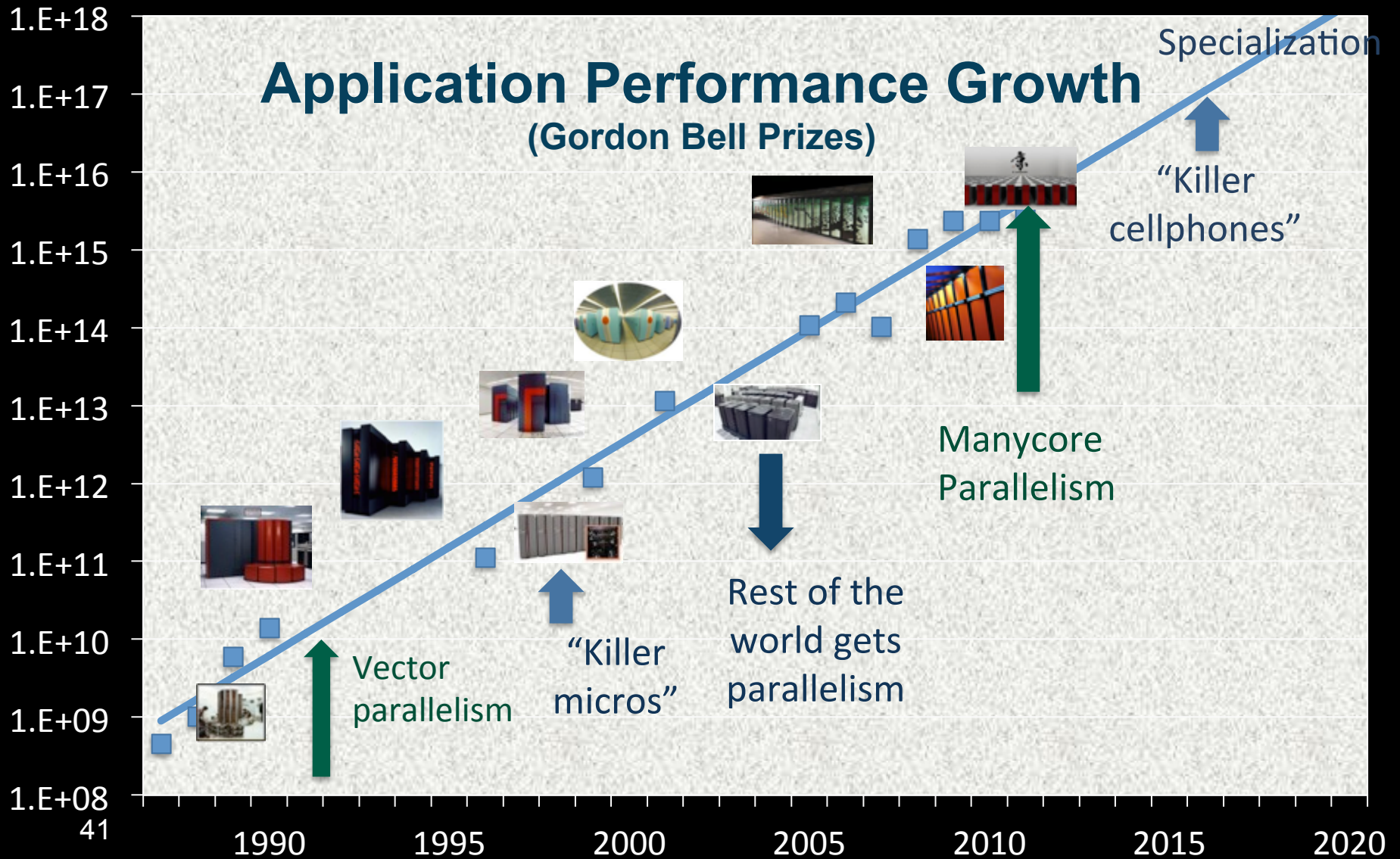
All-to-All



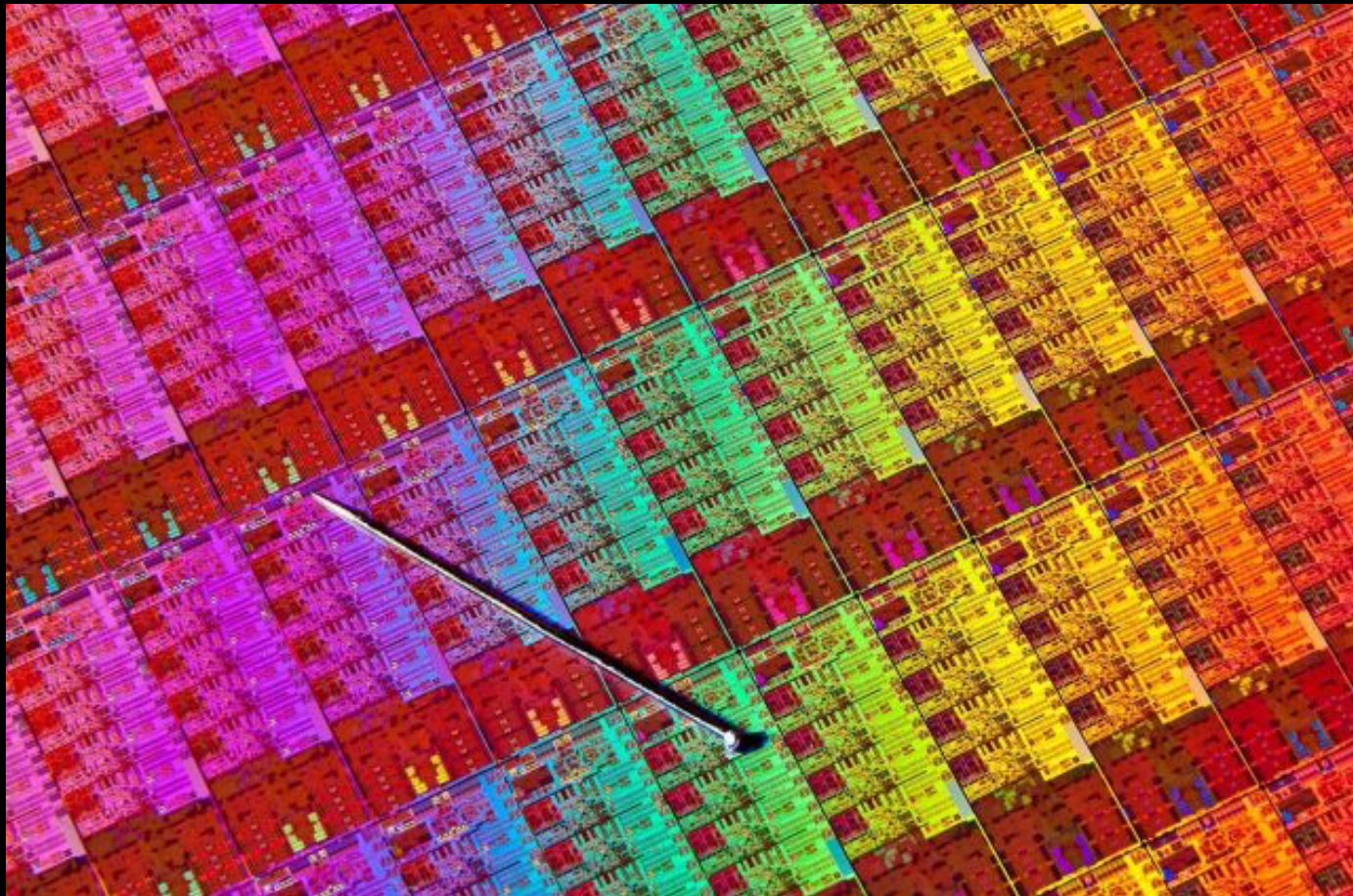
**Random
Access**

**Different architectures for simulation? Can
simulation use data architectures?**

More Parallelism at Lower Levels

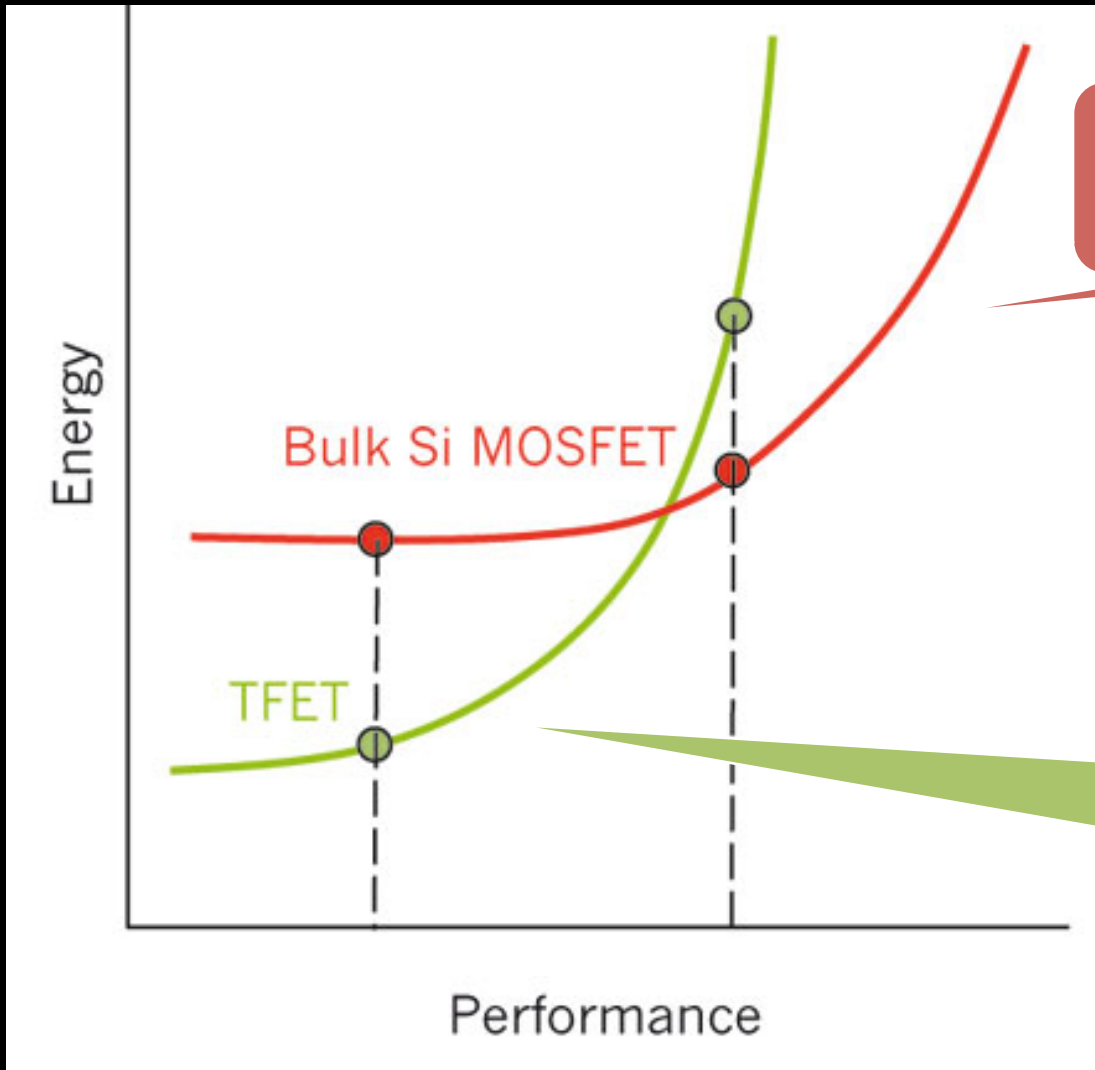


End of Transistor Density Scaling



ITRS now sets the end of transistor shrinking to the year 2021

Device alternatives require lower clock → more parallelism



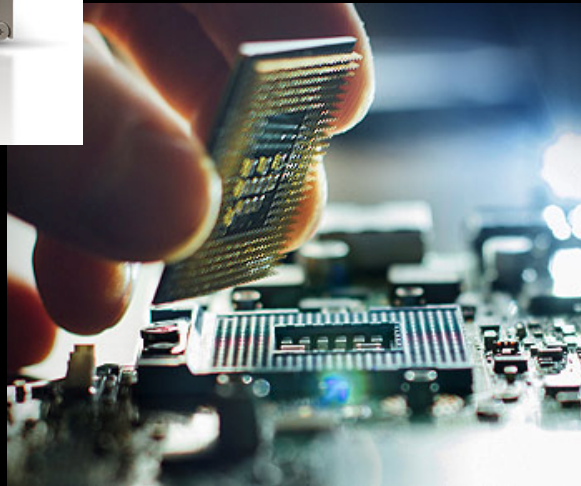
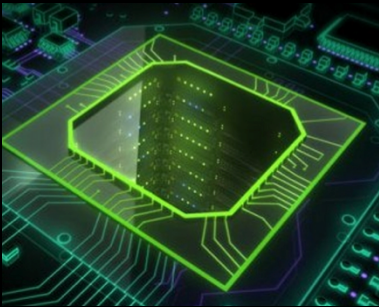
Today's CMOS Technology

Tunneling FET advantage *only at low clock rates*

Specialization: End Game for Moore's Law



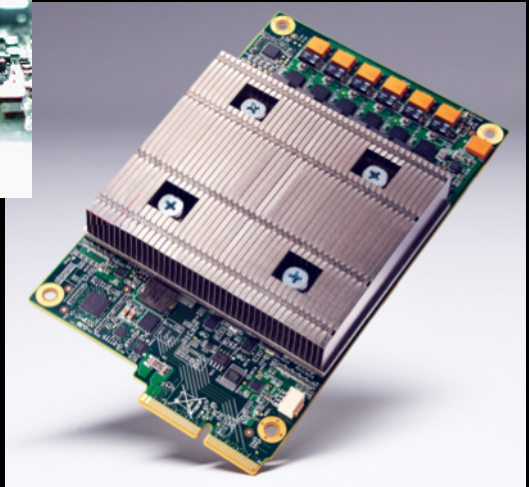
NVIDIA builds deep learning appliance with P100 Tesla's



Intel buys deep learning startup, Nervana



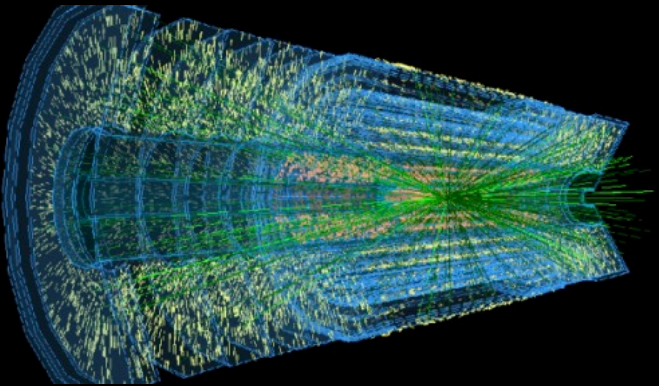
FPGAs



Google designs its own Tensor Processing Unit (TPU)

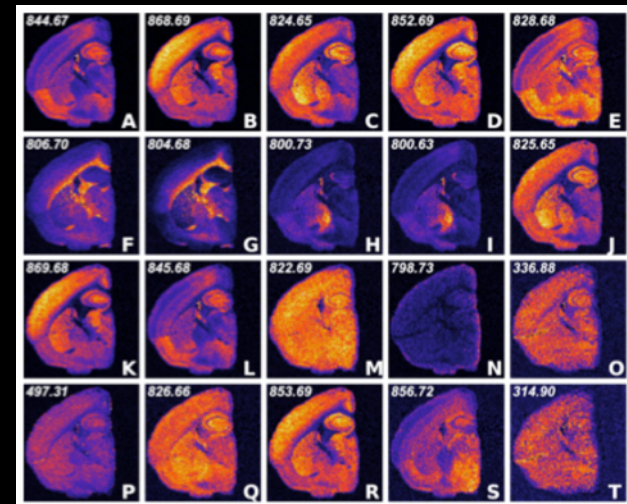
Data processing with special purpose hardware

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



Particle Tracking with Neuromorphic chips

Computing in Detectors



Deep learning processors for image analysis

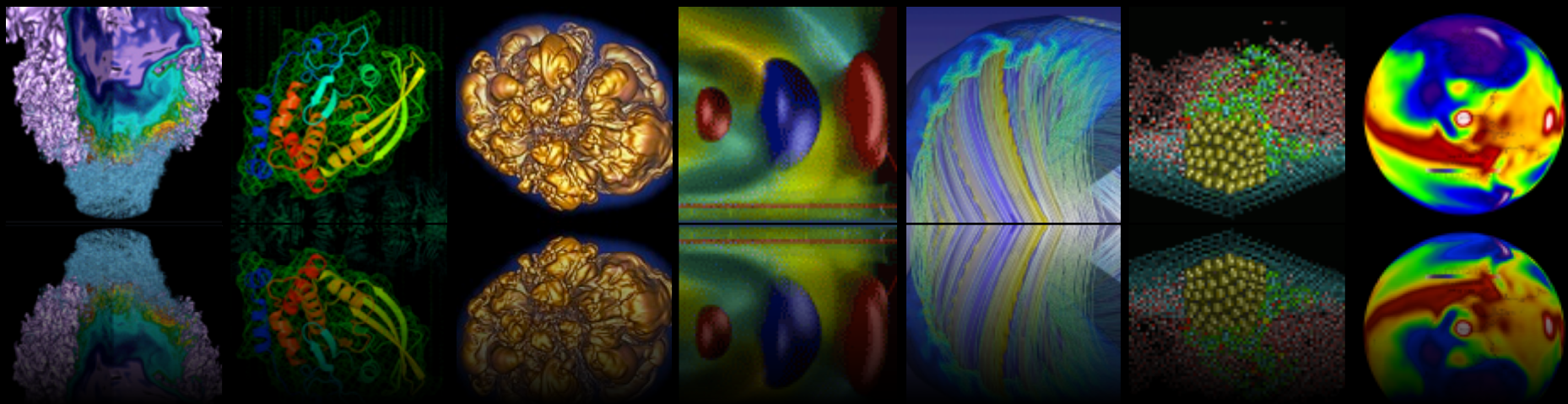
FPGAS for genome analysis

KATHY YELICK'S
2031:
a science odyssey



Life of a Scientist in 2031

- **No personal/departmental computers**
- **Users don't login to HPC Facilities**
- **Travel replaced by telepresence**
- **Lecturers teach millions of students**
- **Theorems proven by online communities**
- **Laboratory work is outsourced**
- **Experimental facilities are used remotely**
- **All scientific data is (eventually) open**
- **Big science and team science democratized**



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.