

Scientific Computing: An Intellectual Lever for Discovery



Dan Reed
reed@renci.org



Chancellor's Eminent Professor
Vice Chancellor for Information Technology
University of North Carolina at Chapel Hill

Director, Renaissance Computing Institute (RENCI)
Duke, UNC Chapel Hill and North Carolina State University

Chair, Board of Directors
Computing Research Association (CRA)



The Instruments of Innovation



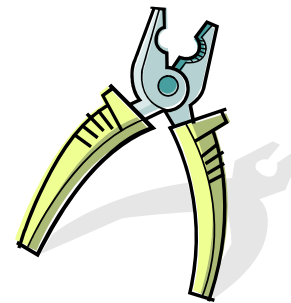
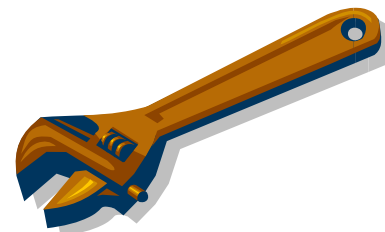
Nothing tends so much to the advancement of knowledge as the application of a new instrument. The native intellectual powers of men in different times are not so much the causes of the different success of their labors, as the peculiar nature of the means and artificial resources in their possession.

Sir Humphrey Davy

The Computer Science Behind Science

- **Computer science (CS) enables science**

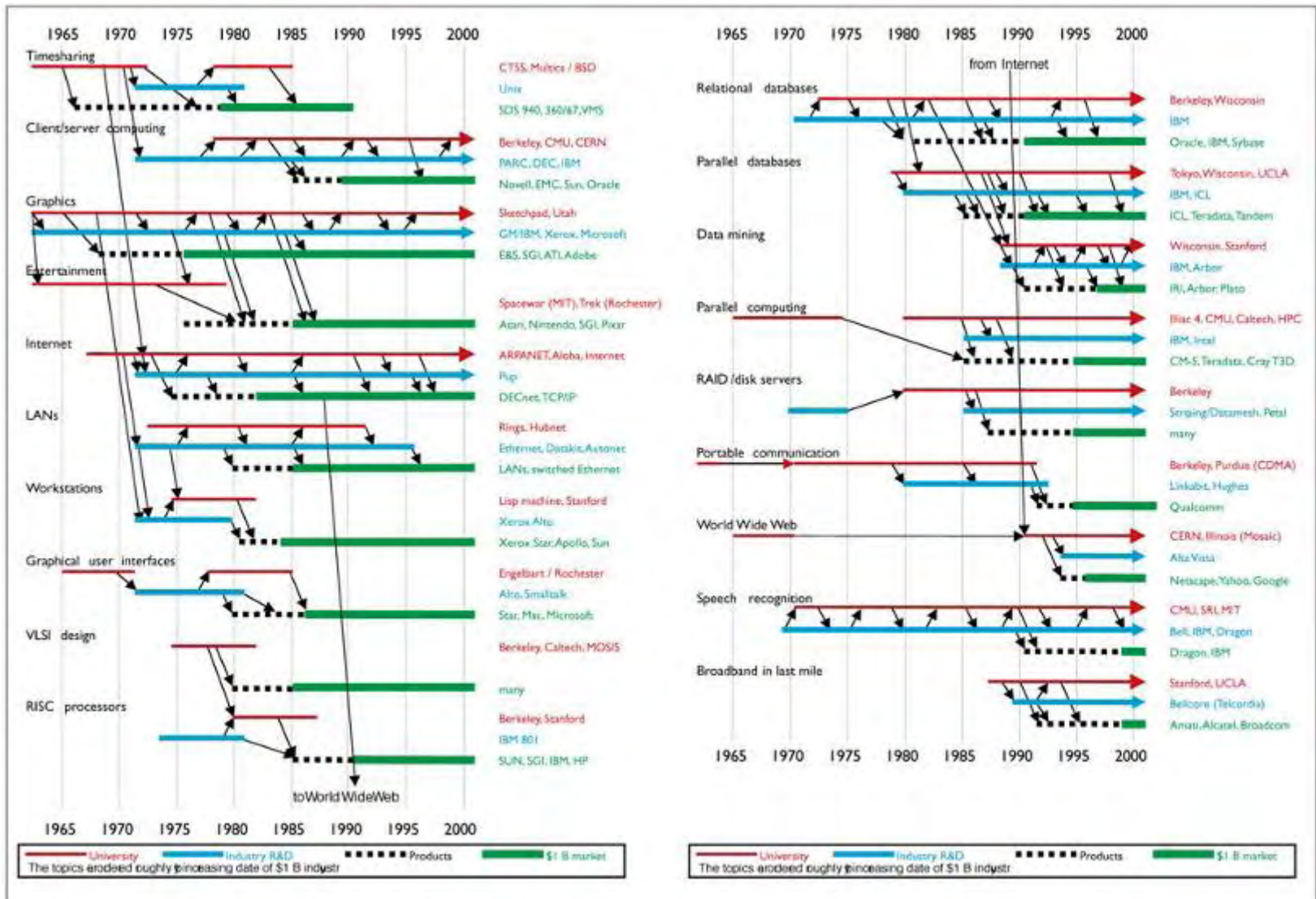
- high-performance computing
- networks and sensors
- grids and web services
- data models and mining
- scientific and information visualization
- distributed collaboration tools
- algorithms, software and tools
- numerical analysis
- artificial intelligence



- **Let's look at some examples ...**

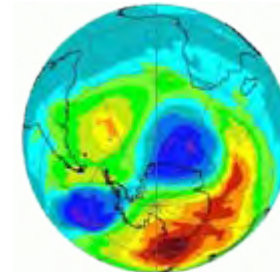
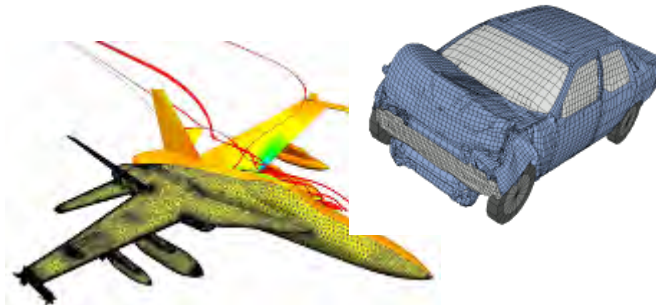
- where CS and science meet as computational science

Computer Science Impact



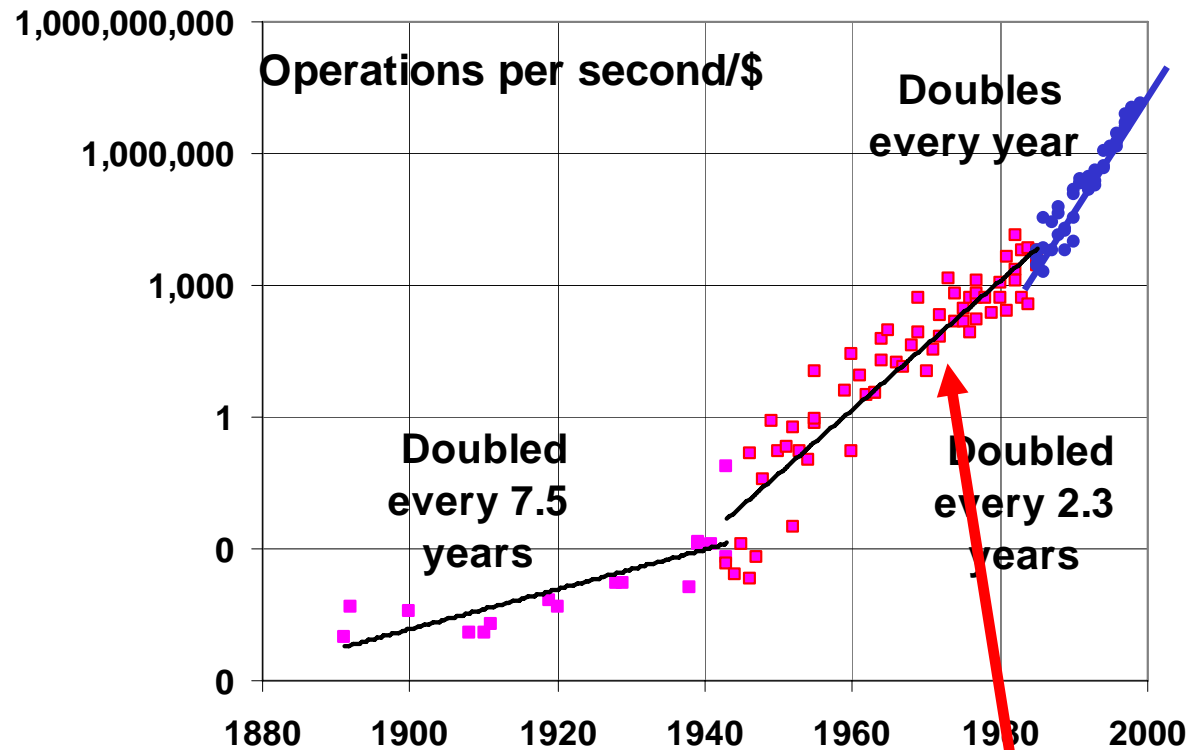
The Third Pillar of 21st Century Science

- **Three pillars**
 - theory, experiment and *computational science*
- **Computational science enables us to**
 - investigate phenomena where
 - economics or constraints preclude experimentation
 - evaluate complex models and manage massive data model processes across interdisciplinary boundaries
 - transform business and engineering practices

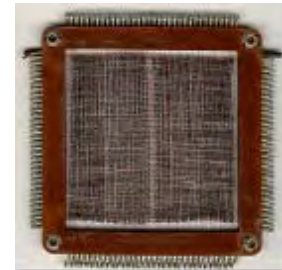


Computing History and Exponentials

- **1890-1945**
 - mechanical, relay
 - 7 year doubling
- **1945-1985**
 - tube, transistor,..
 - 2.3 year doubling
- **1985-2003**
 - microprocessor
 - 1 year doubling



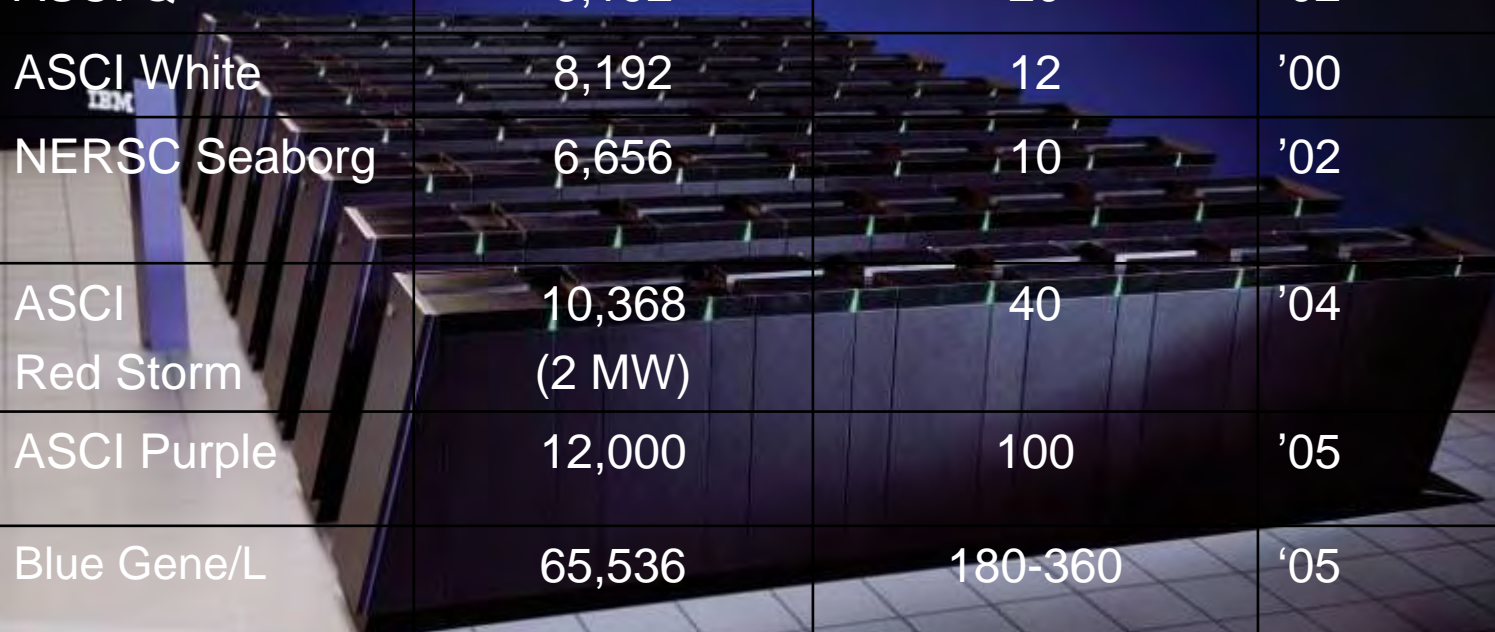
- **Exponentials**
 - chip transistor density: 2X in ~18 months
 - graphics: 100X in three years
 - WAN bandwidth: 64X in two years
 - storage: 7X in two years



4K bit core plane

Microprocessor Revolution

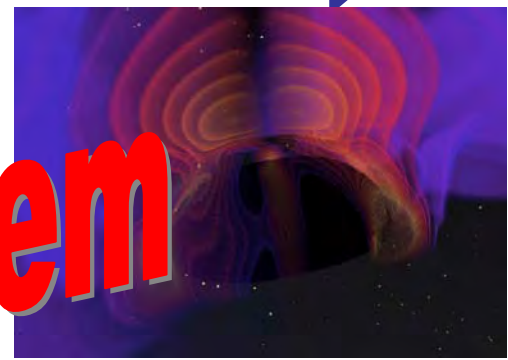
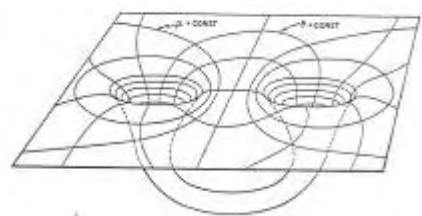
Large Systems



Machine	Processor Count	Teraflops	Year
ASCI Q	8,192	20	'02
ASCI White	8,192	12	'00
NERSC Seaborg	6,656	10	'02
ASCI Red Storm	10,368 (2 MW)	40	'04
ASCI Purple	12,000	100	'05
Blue Gene/L	65,536	180-360	'05
NASA Columbia	10,240 (2 MW)	60	'04

Black Hole Collision Problem

1,800,000,000X



A "Clean" Problem

1963
Hahn and Lindquist
IBM 7090
One Processor
Each 0.2 MF
3 Hours

1977
Epple and Smarr
CDC 7600
One Processor
Each 35 MF
5 Hours

1999
Seidel and Suen, et al.
NCSA SGI Origin
256 Processors
Each 500 MF
40 Hours

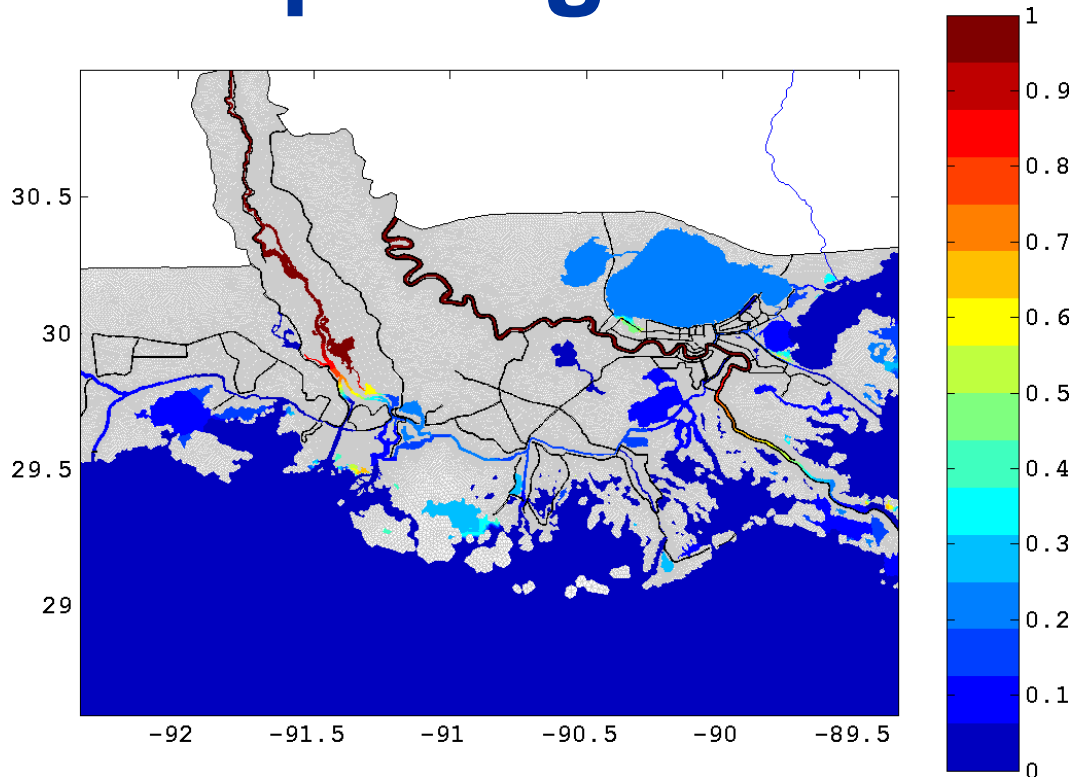
2001
Seidel et al
NCSA Pentium III
256 Processors
Each 1 GF
500,000 Hours total
plus 500,000 hours at NERSC

300X

30,000X


~200X

Computing for Disaster Response

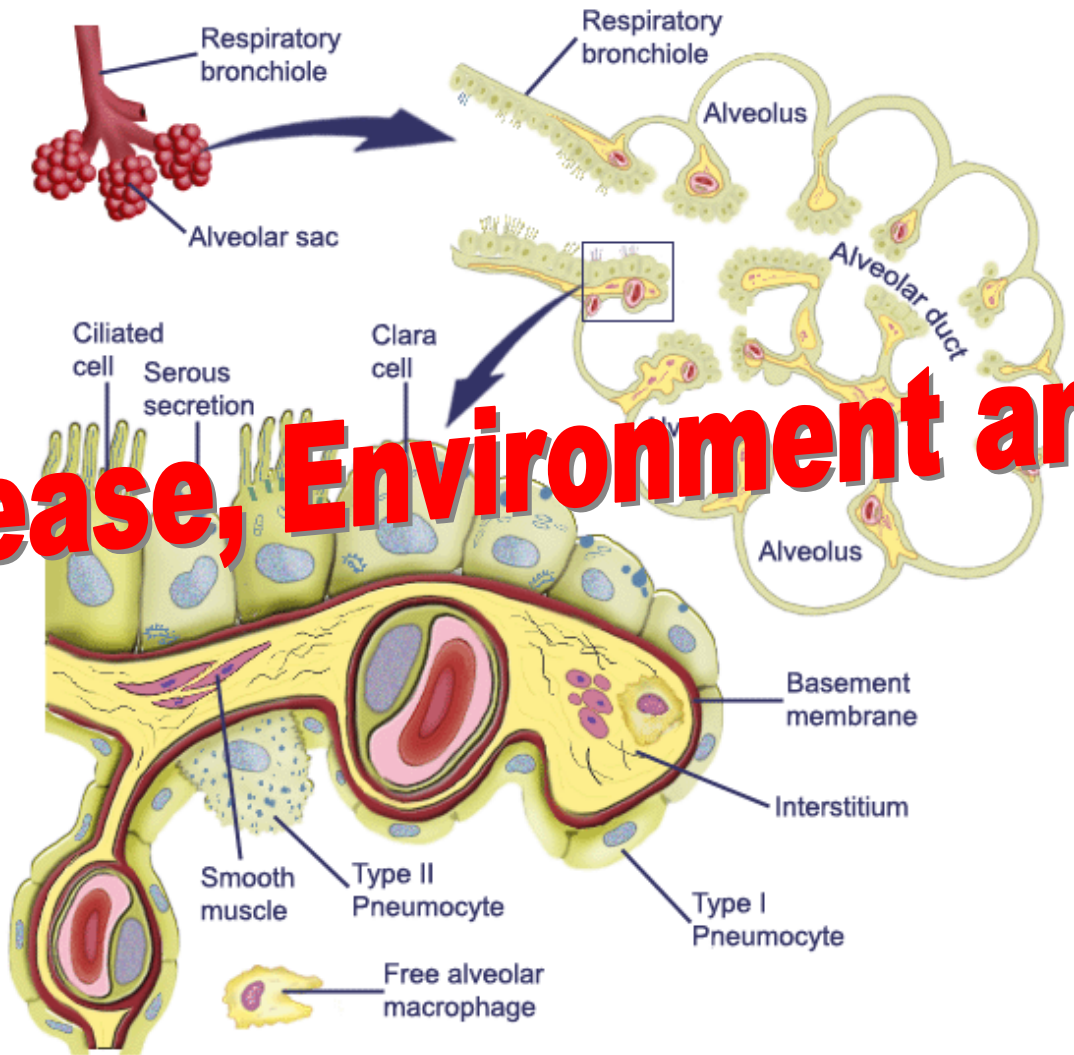


- **Post-Katrina NOAA challenge**
 - petrochemical spillage and remediation
 - water levels determine dissemination
 - answers needed within 48 hours, but inadequate computing capability

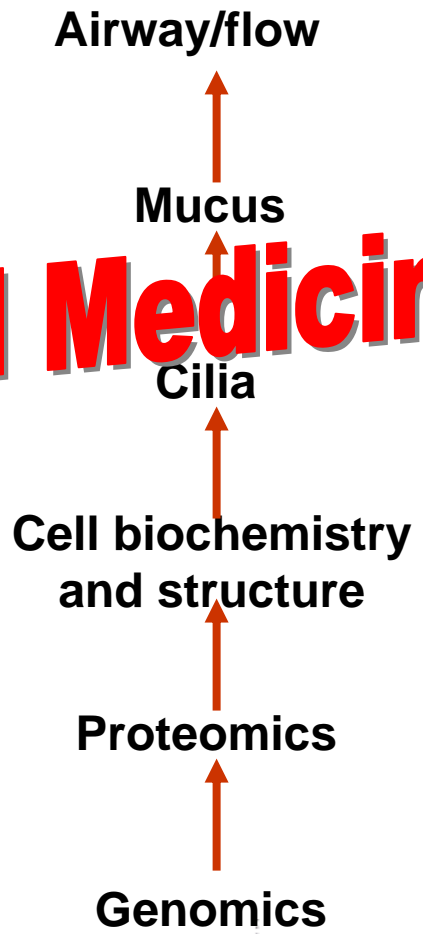
- **UNC Marine Sciences and RENCi**

 ADCIRC storm surge model predicts water levels
HPC system yields model predictions

Biophysical and Environmental Modeling



Disease, Environment and Medicine



Source: Ric Boucher, UNC



The Coming of Massive Parallelism

- **Technology trends**

- multicore processors

- IBM Power5/6 and SUN UltraSPARC IV
 - Intel Xeon and AMD Opteron
 - quad core and beyond are coming

- reduced power consumption

- laptop and mobile market drivers

- greater I/O and memory integration

- PCI Express, Infiniband, ...

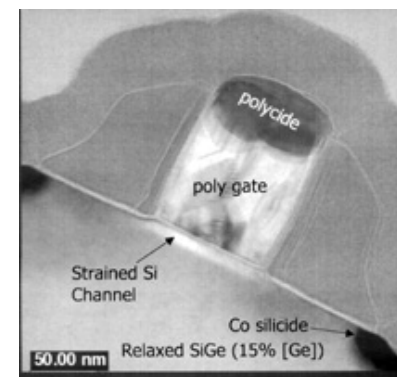
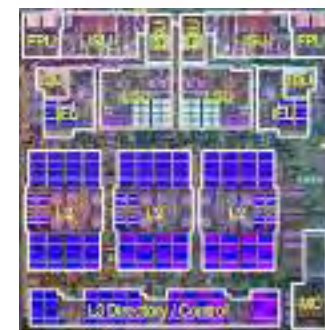
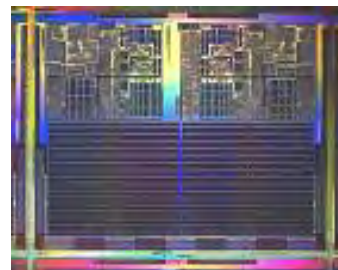
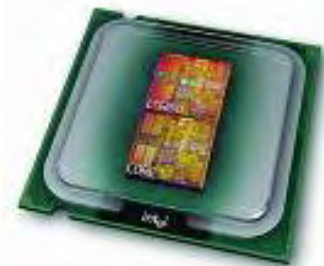
- Justin Ratter (Intel)

- “100’s of cores on a chip in 2015”

- **Moore’s law isn’t a birthright**

- CMOS scaling issues are now a challenge

- power, junction size, fab line costs, ...



Digital Reality: The Exponentials



1956



1972

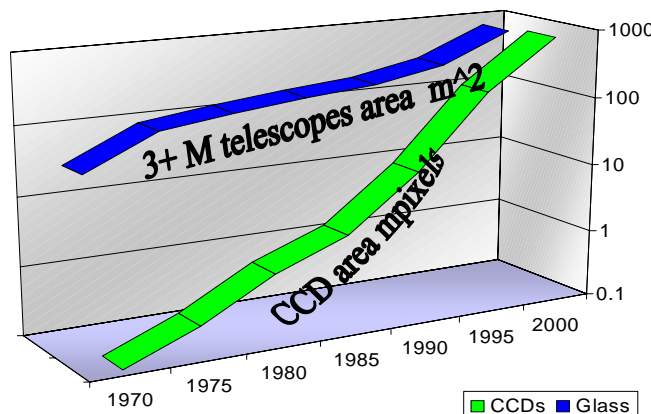
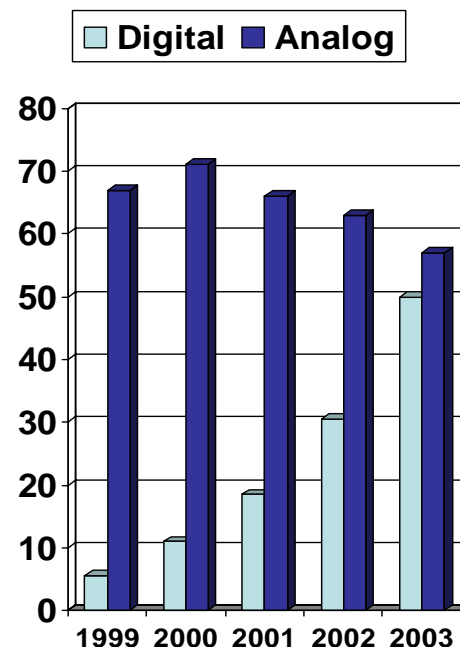


2006

- **Megabyte**
 - a small novel
- **Gigabyte**
 - a pickup truck filled with paper or a DVD
- **Terabyte: one thousand gigabytes** – ~\$1000 today
 - the text in one million books
 - entire U.S. Library of Congress is ~ten terabytes of text
- **Petabyte: one thousand terabytes**
 - 1-2 petabytes equals all academic research library holdings
 - coming soon to a pocket near you!
 - *soon routinely generated annually by many scientific instruments*
- **Exabyte: one thousand petabytes**
 - 5 exabytes of words spoken in the history of humanity

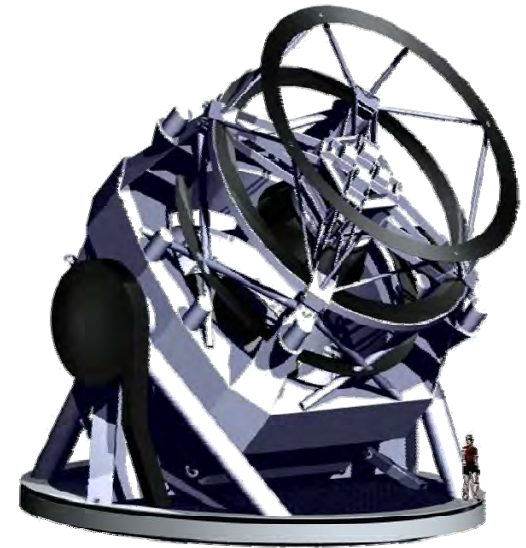
Astronomy and Consumer Cameras

- **Digital camera sales**
 - now exceed analog
 - January 2006
 - Nikon stops film camera production
- **From glass plates to CCDs**
 - detectors follow Moore's law
 - data tsunami
 - data doubles every two years
- **Telescope growth**
 - 30X glass (concentration)
 - 3000X in pixels (resolution)
- **Single astronomy images**
 - 16Kx16K pixels and growing
- **Detector driver**
 - consumer electronics



Large Synoptic Survey Telescope (LSST)

- **Top project of the astronomy decadal survey**
- **Celestial cinematography**
 - 3 gigapixel detector for wide field imaging
- **Science**
 - beyond the standard model
 - non-baryonic dark matter
 - non-zero Λ and neutrino oscillations
 - observation targets
 - near Earth object survey
 - weak lensing of wide fields
 - supernovae measurements
- **Features**
 - 9.6 square degree field/6.5 meter effective aperture
 - *~15 TB of data/night, target first light 2012*

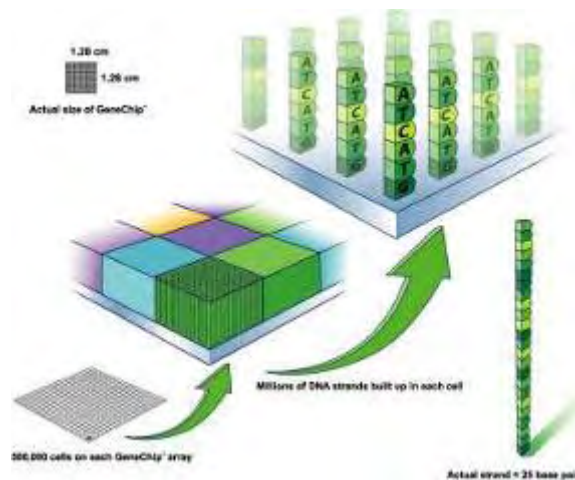


Gene Expression and Microarrays

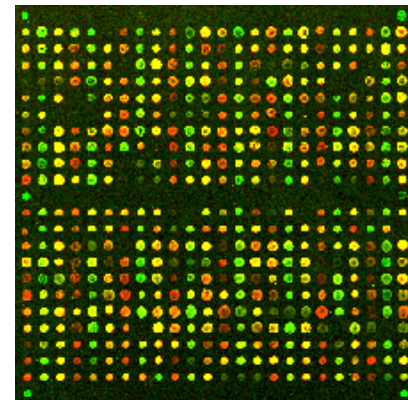
- **Concurrent evaluation**
 - expression levels for thousands of genes
- **Photolithography**
 - up to 500K 10-20 micron cells
 - each containing millions of identical DNA molecules
- **Image capture and analysis**
 - laser scanning and intensity calculation



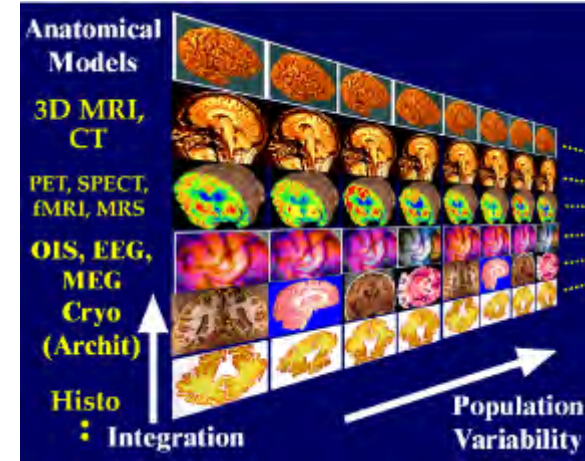
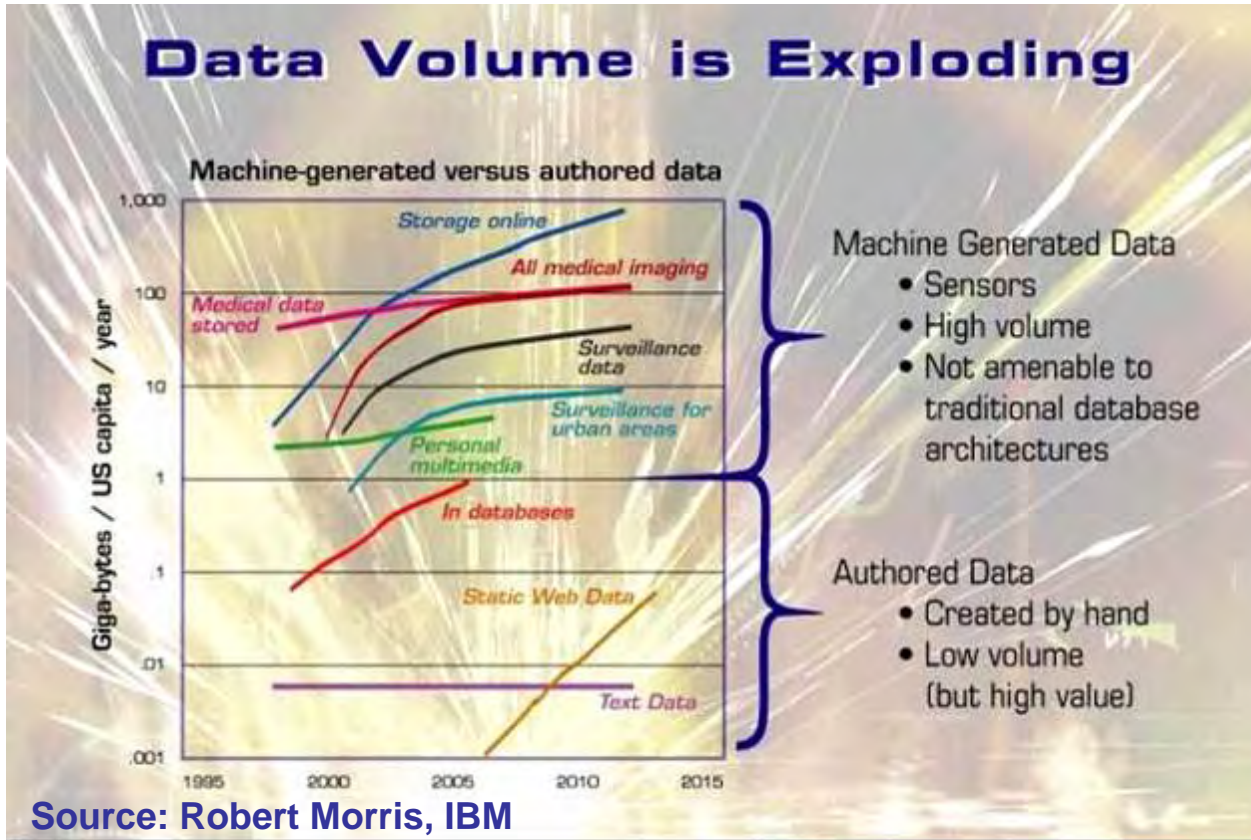
AFFYMETRIX GENECHIP®



Source: Affymetrix



Sensor Data Overload



Source: Chris Johnson, Utah
Art Toga, UCLA

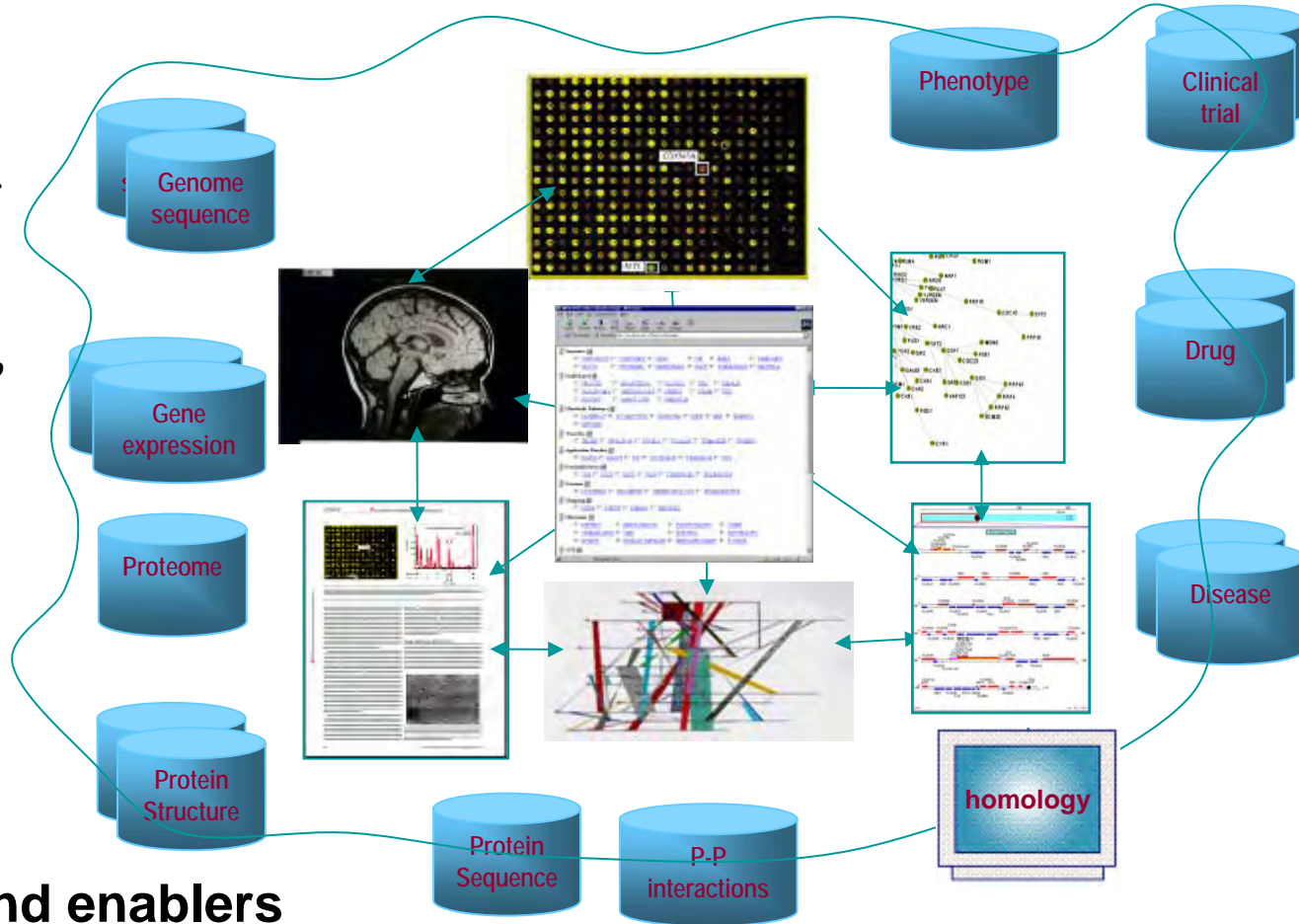
- High resolution brain imaging

4.5 petabytes (PB) per brain



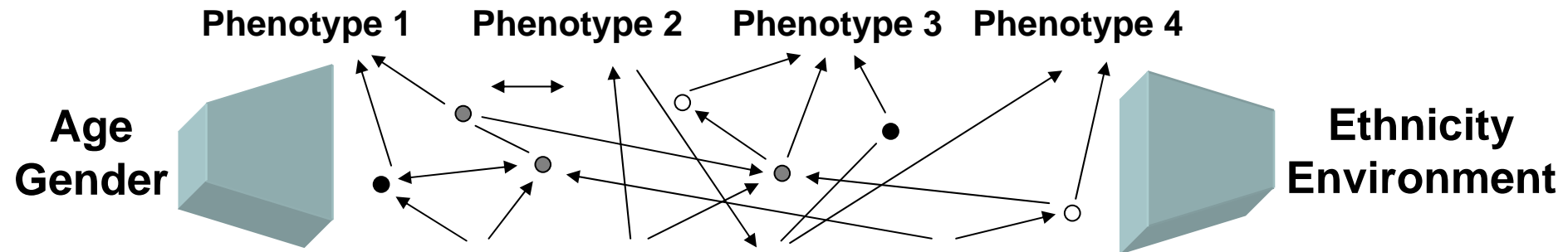
Data Heterogeneity and Complexity

Genomic, proteomic, transcriptomic, metabolomic, protein-protein interactions, regulatory bio-networks, alignments, disease, patterns and motifs, protein structure, protein classifications, specialist proteins (enzymes, receptors)

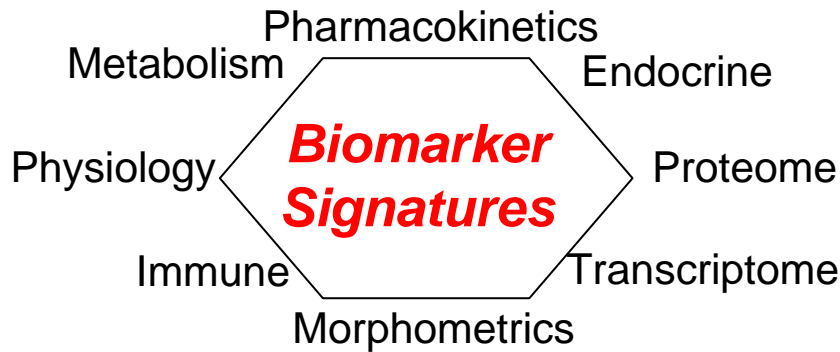


- Many causes and enablers
 - increased instrument resolution
 - increased storage capability

Genetics and Advanced Data Mining



Identify Genes



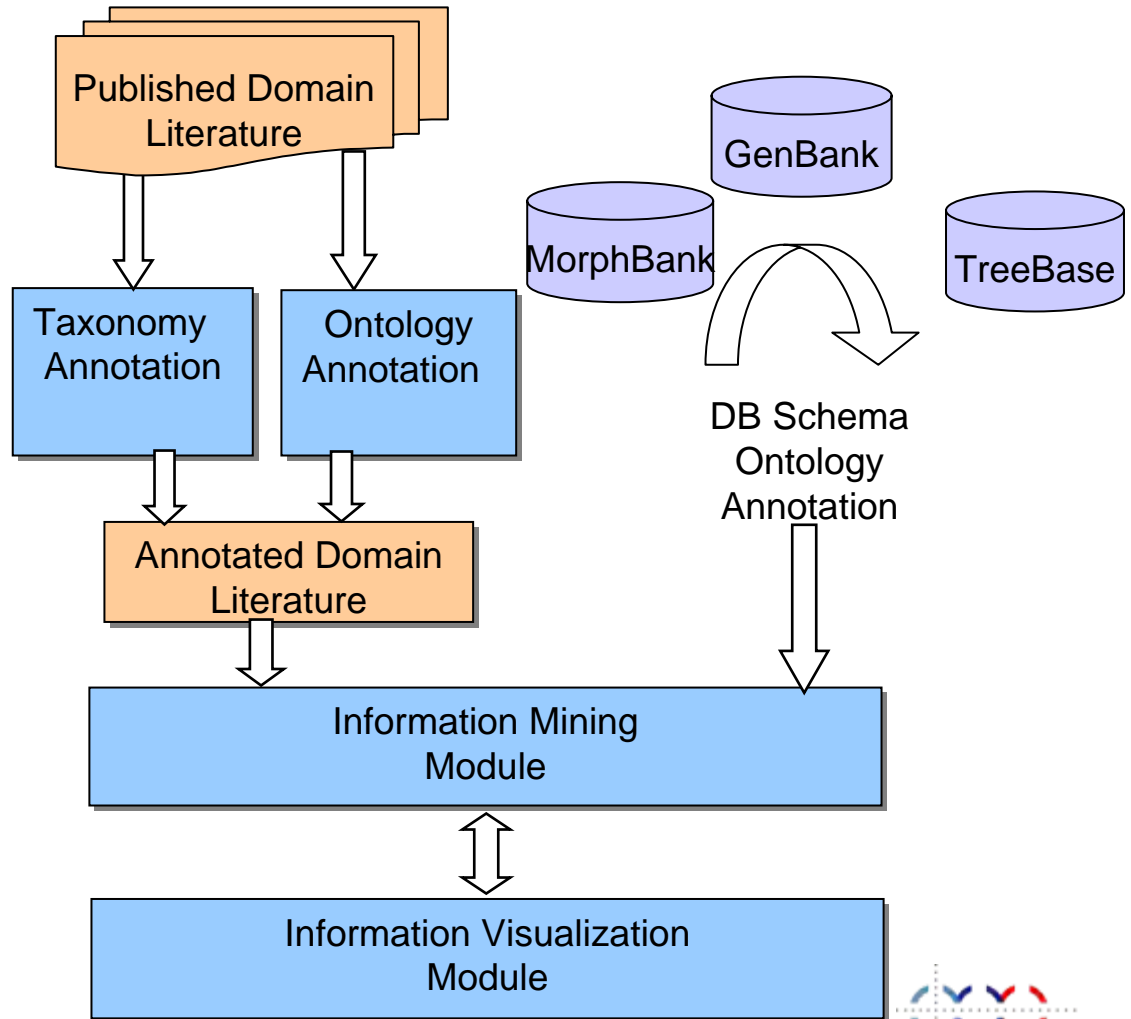
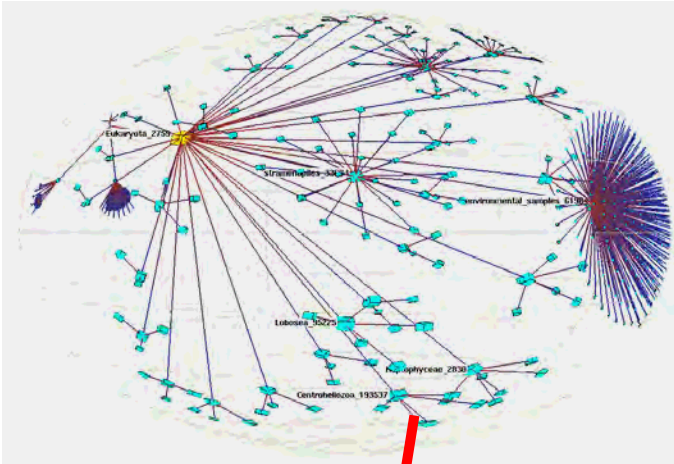
Predictive Disease Susceptibility



Source: David Threadgill/Terry Magnuson, UNC

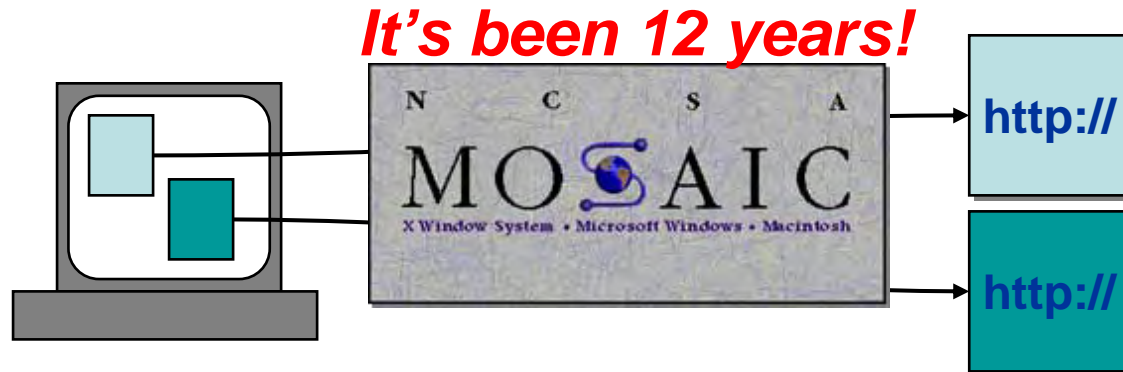


Data Federation and Info Viz

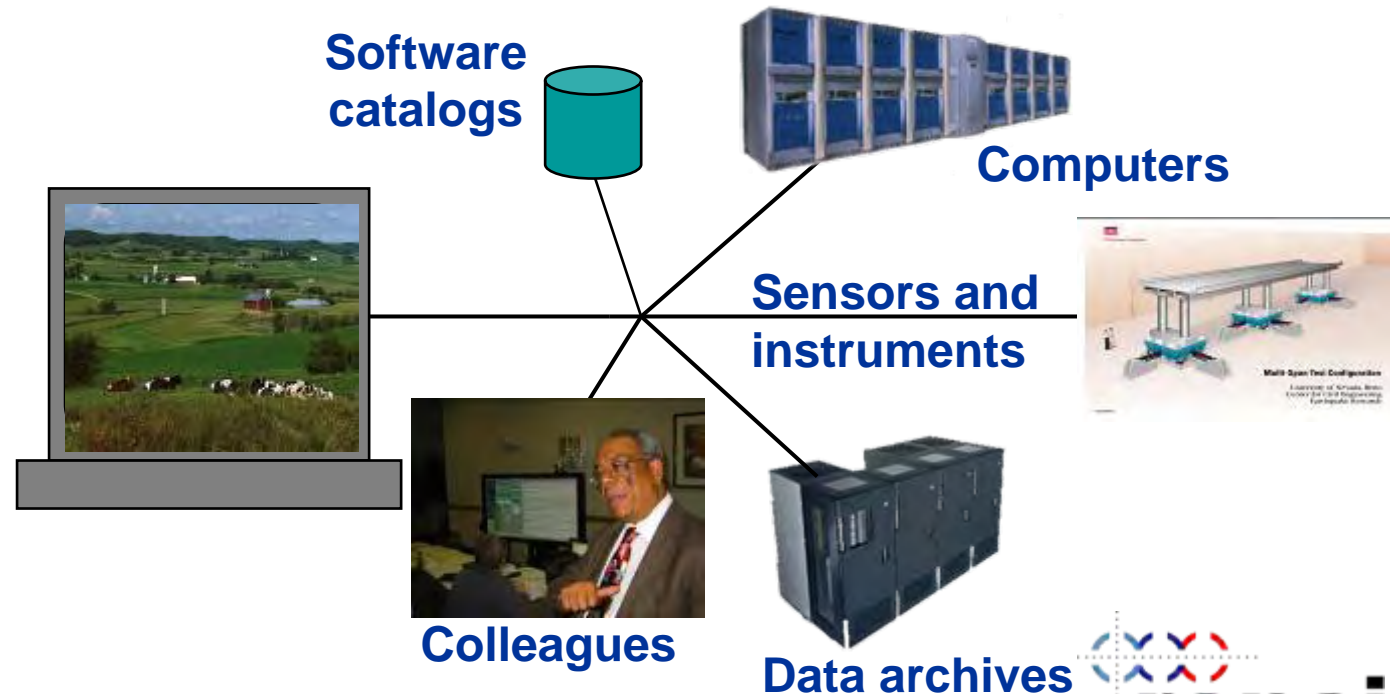


What's A Grid?

Web: Uniform access to documents

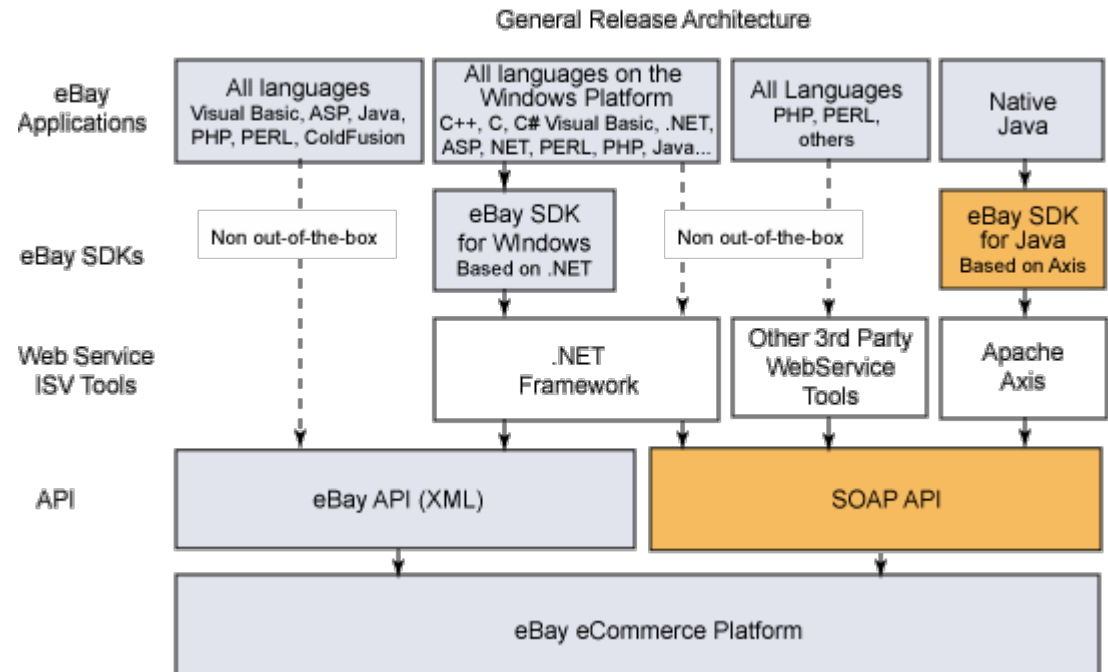
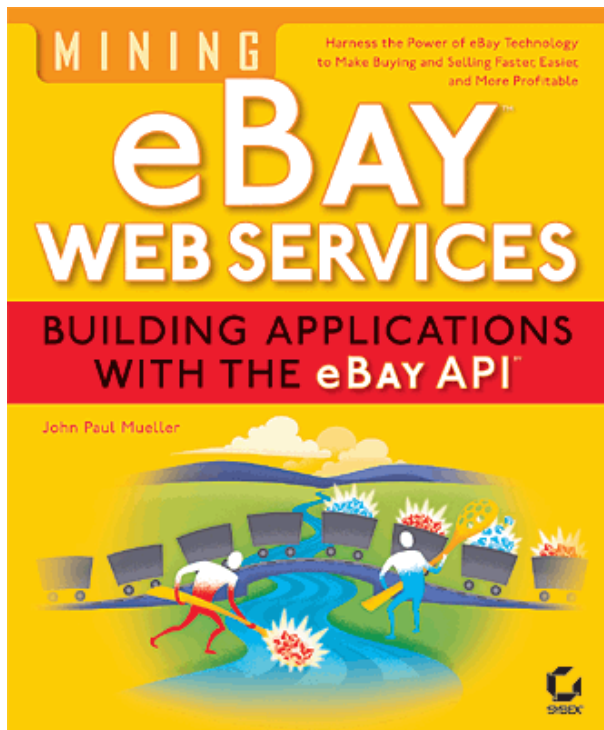


Grid: Flexible, high-performance access to resources and services for *distributed communities*



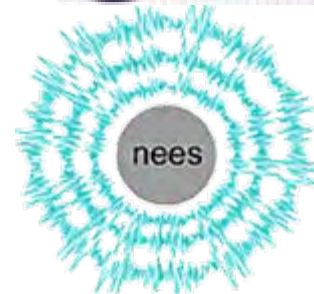
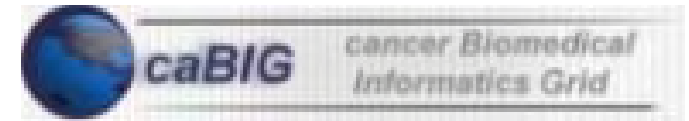
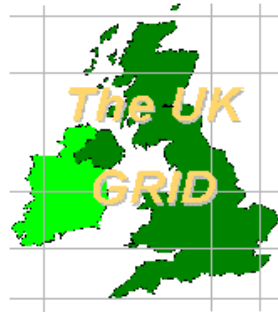
eBay Web Services Architecture

- Over 40% of eBay's listings are now via API calls

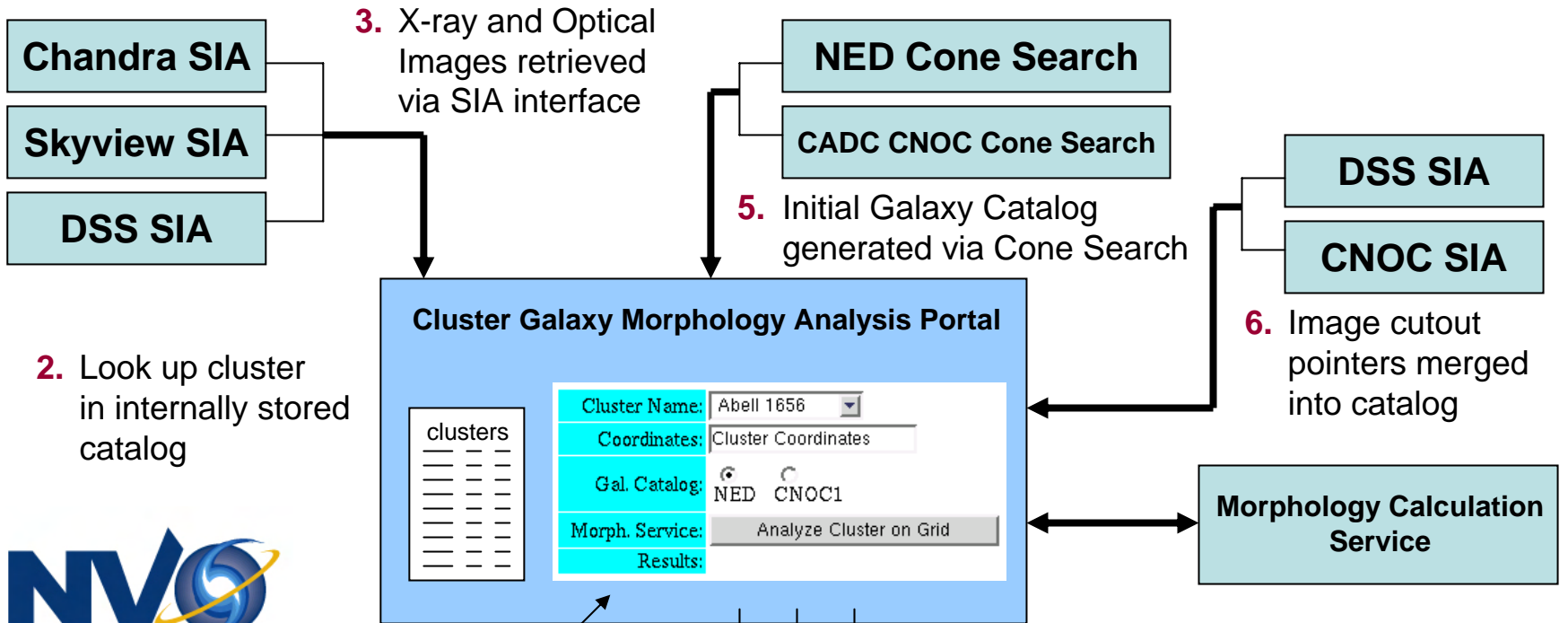


Source: IBM

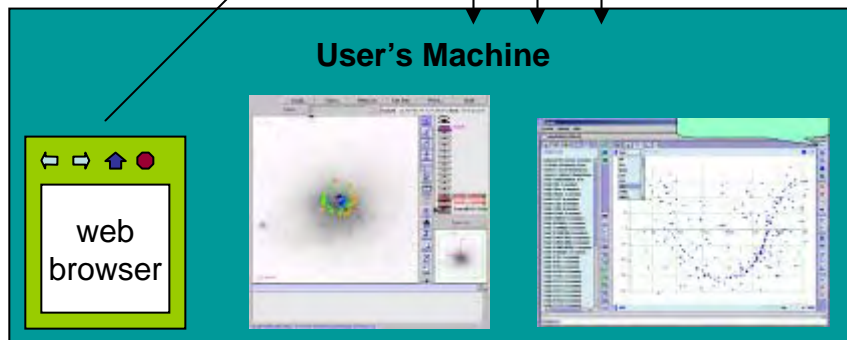
Grids/Web Services of All Flavors



{Inter}national Virtual Observatory



1. User selects a cluster
2. Look up cluster in internally stored catalog
3. X-ray and Optical Images retrieved via SIA interface
4. User launches distributed analysis
5. Initial Galaxy Catalog generated via Cone Search
6. Image cutout pointers merged into catalog
7. Morphological parameters calculated on grid for each galaxy
8. User downloads final table and images for analysis & visualization

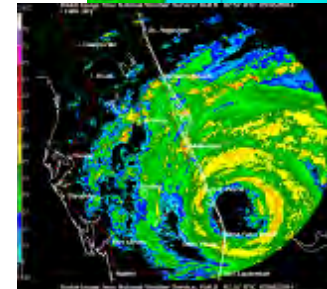


Source: Ray Plante, NCSA

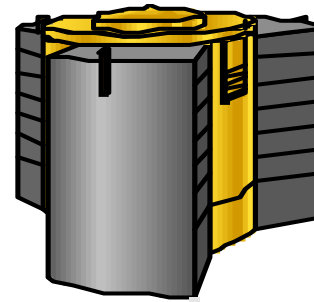
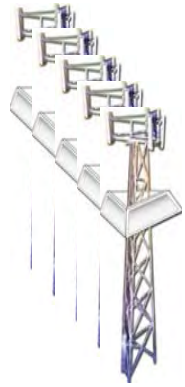


Weather and Economic Loss

- **\$10T U.S. economy**
 - 40% is adversely affected by weather and climate
- **\$1M in loss to evacuate each mile of coastline**
 - we now over warn by 3X!
 - average over warning
 - 200 miles, or \$200M per event
- **Improved forecasts**
 - lives saved and reduced cost
- **LEAD national Grid**
 - Oklahoma, Indiana, UCAR
 - Colorado State, Howard, Alabama
 - Millersville, NCSA, North Carolina



The LEAD Vision: A Paradigm Shift

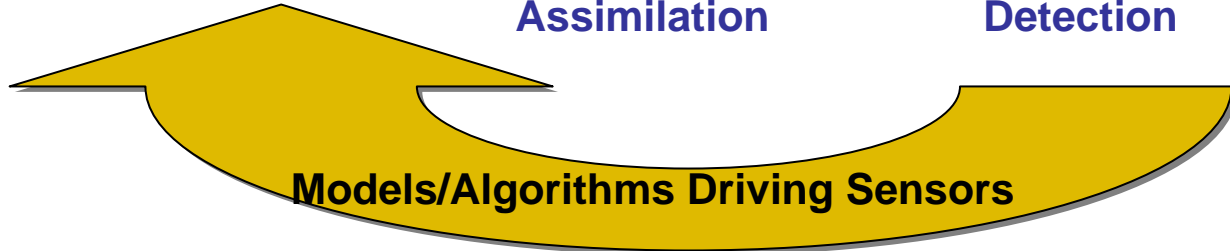


Dynamic Observations

Analysis and Assimilation

Prediction and Detection

Product Generation and Dissemination

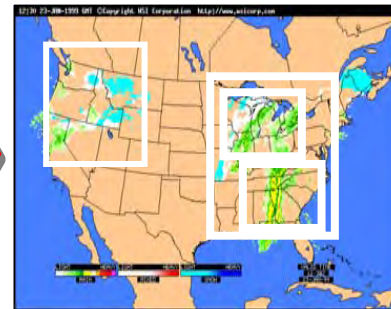
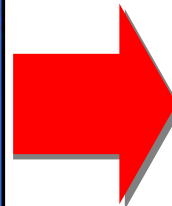
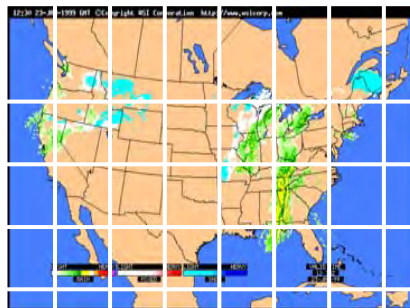


Models/Algorithms Driving Sensors

The challenge: Build cyberinfrastructure services that provide adaptability, scalability, availability, usability, and real-time response.



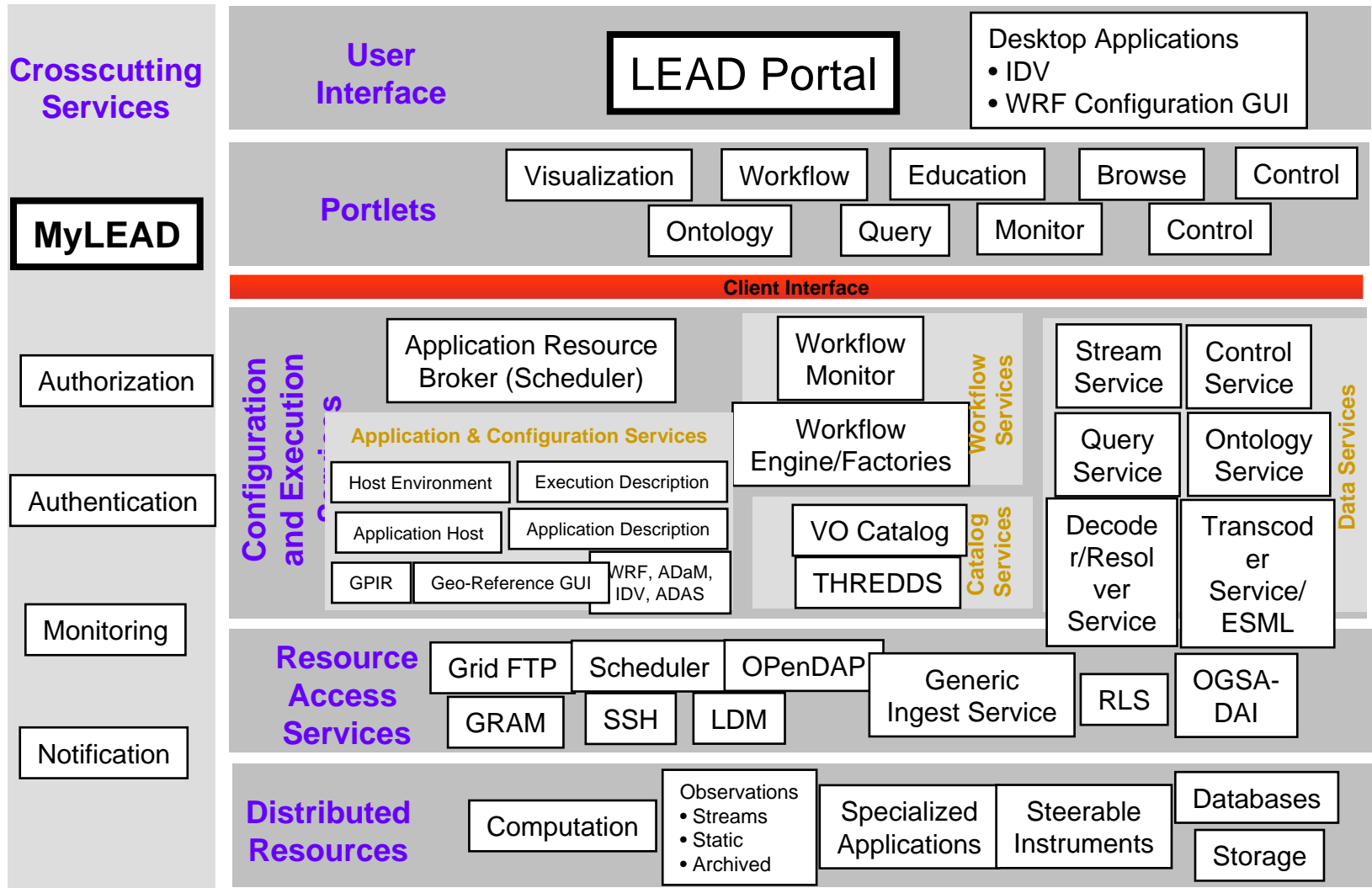
End Users



Source: Plale, Indiana



LEAD Service Oriented Architecture



LEAD Mesoscale Meteorology

LEAD PORTAL
Linked Environments for Atmospheric Discovery
Sponsored by the National Science Foundation

Portal Home | Geo GUI | Education and Outreach | Weather | Links | About LEAD | Help

Home

To view a local radar, select area of interest and click on the image below.
**RADAR REFLECTIVITY FROM RADAR CODED MESSAGES
 NATIONAL WEATHER SERVICE
 AUTOMATED EDITING APPLIED
 SEP 24, 2005 21:49 UTC**

MEG
 55 DBZ
 50 DBZ
 45 DBZ
 40 DBZ
 30 DBZ
 15 DBZ

Data provided by NOAA's National Weather Service

LEAD Home | FAQ | Privacy | Terms of use | Contact us

User Name:

Password:

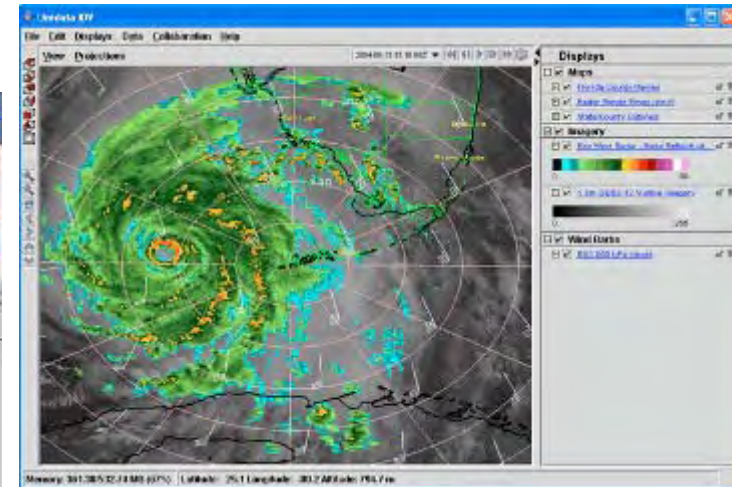
Remember my login

[Create new account](#)
[Forgot your password?](#)

LEAD Grid Testbed Status

Testbed	Grid	Auth	GRAM	GridFTP
IU [chinkapin]	✓	✓	✓	✓
NCSA [copper]	✓	✓	✓	✓
OU [aquaman]	✓	✓	✓	✓
GAH [frozone]	✓	✓	✓	✓
UNC [dante0]	✗	✗	✗	✗
Unidata [lead1]	✓	✓	✓	✓

Last updated: Sat Sep 24 17:00:00 2005 Indiana local



Workflow Composer

Workflow: MyLead | Component: Monitor

Component List

- System Components
- http://www.xdrone.indiana.edu
- adder
- Multiplier
- Divider
- discconn
- forkjob
- seps-frm
- seps-etc
- ev2000-4c

Component Information

Service: decoder

Description:
A service for decoding raw data to netcdf format.

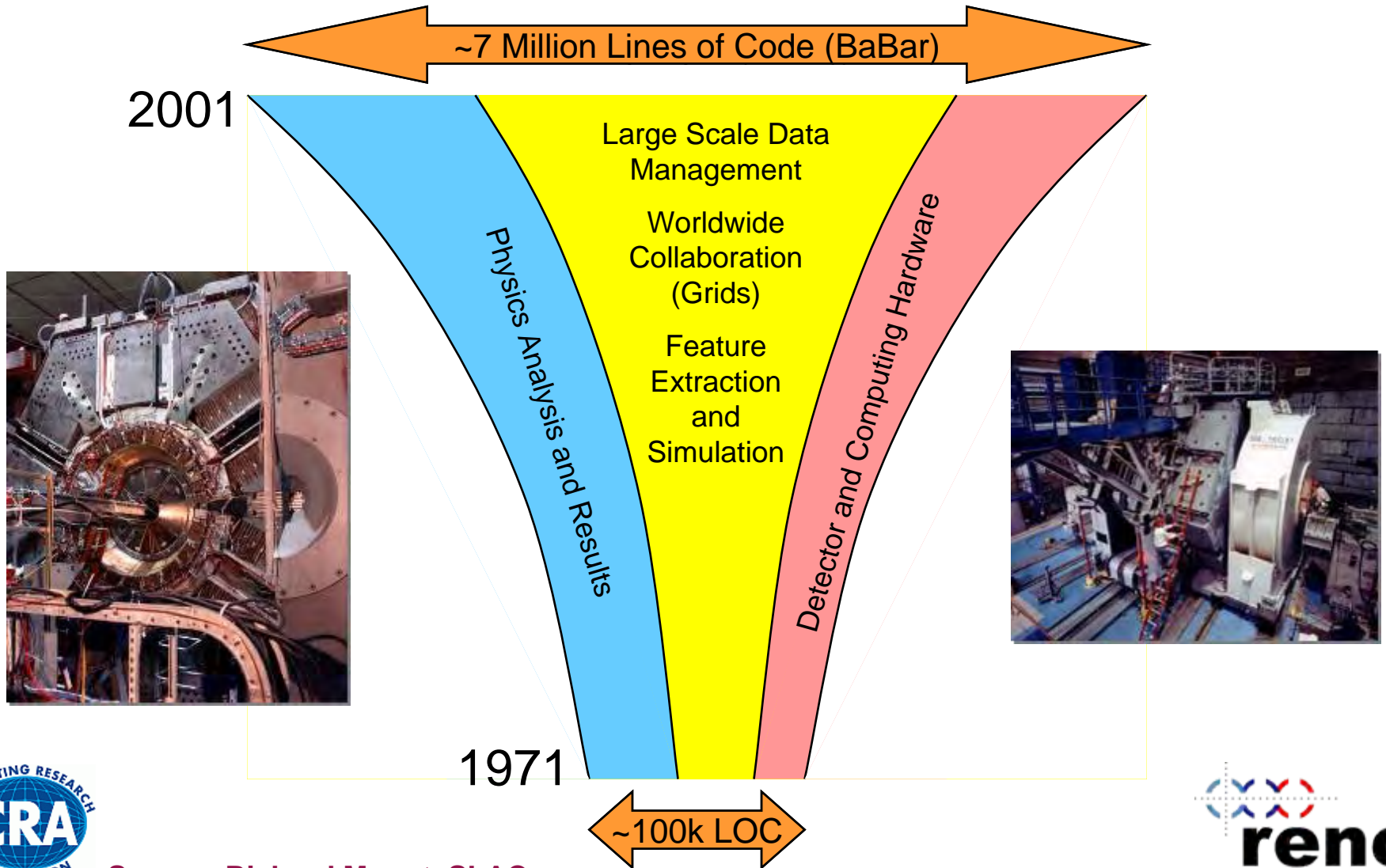
Operation: Run

Selected Output Port

Selected Input Port

Component: Output_URL
 Port: Parameter
 Type: Any
 Description: This port can be connected to any type.

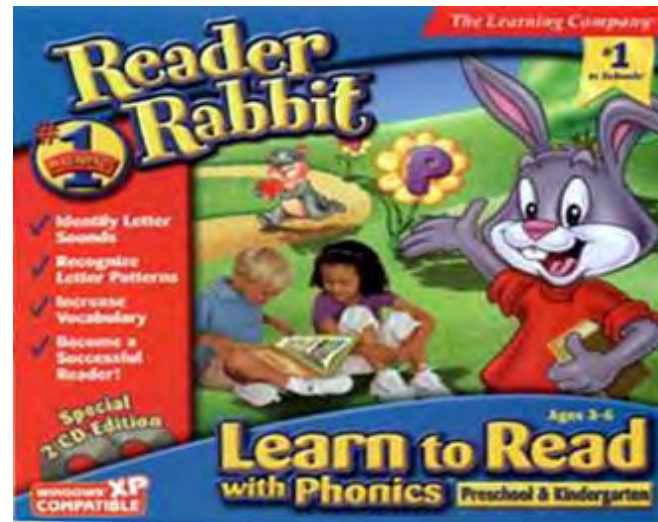
Software Complexity and Collaboration



Need: Simple, Easy-To-Use Tools

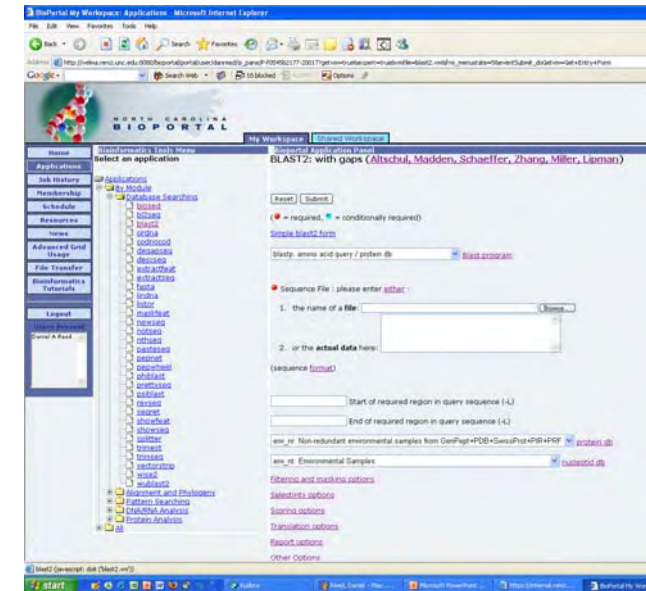
“Genome. Bought the book. Hard to read.”

Eric Lander



Carolina Bioportal

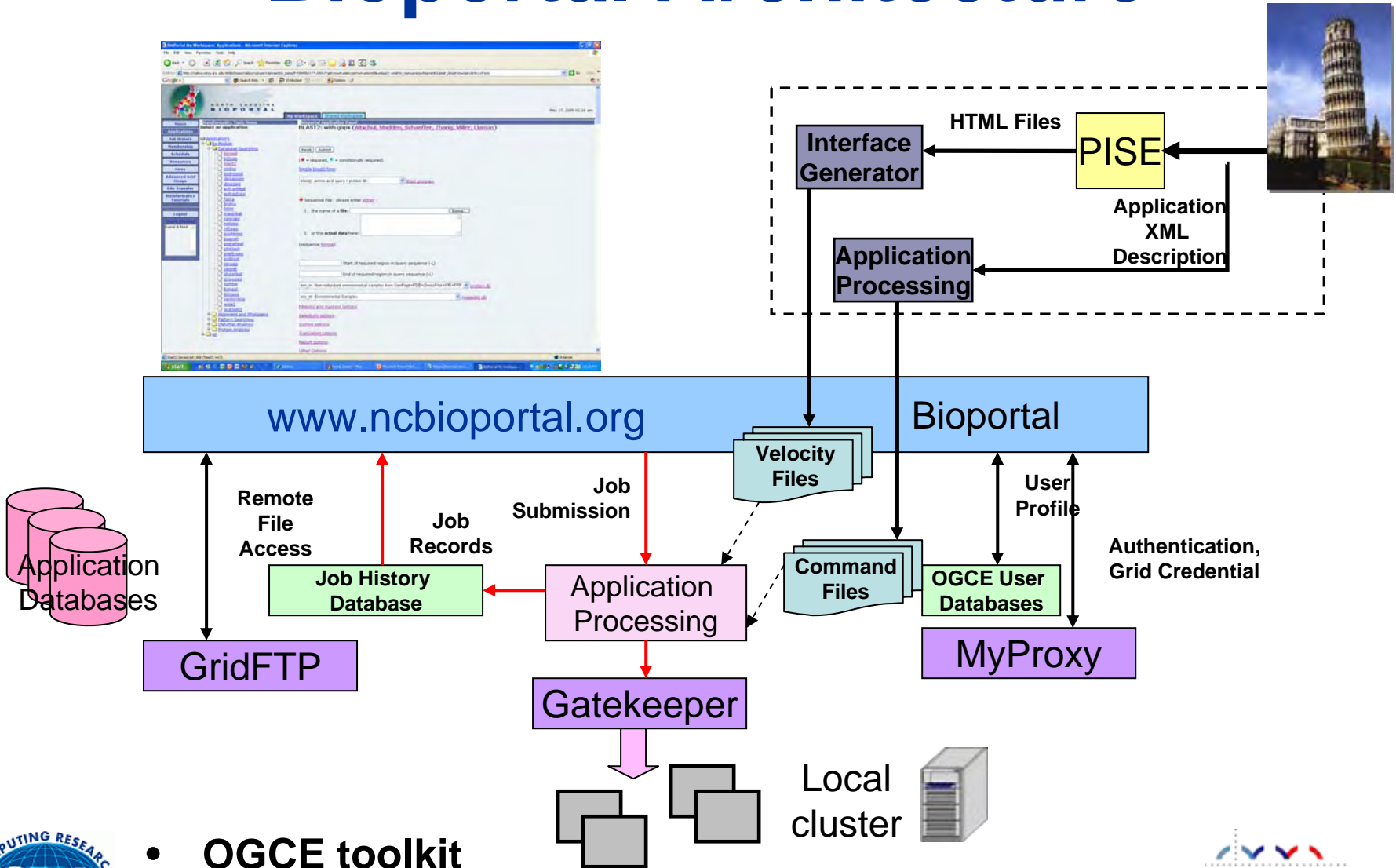
- **Three overlapping target groups**
 - undergraduate education
 - graduate education and research
 - academic/industrial research
- **Features**
 - access to common bioinformatics tools
 - extensible toolkit and infrastructure
 - OGCE and National Middleware Initiative (NMI)
 - leverages emerging international standards
 - remotely accessible or locally deployable
 - packaged and distributed with documentation
- **National reach and community**
 - NSF TeraGrid deployment
 - science gateway
- **Education and training**
 - hands-on workshops
 - clusters, Grids, portals and bioinformatics



North Carolina



Biportal Architecture



- **OGCE toolkit**
 - used by cyberinfrastructure projects
 - LEAD, NEES, PACI, DOE, TeraGrid ...

Information and Social Processes

- **Google**
 - it's a search engine, it's a verb, ...
- **Blogs**
 - published self-expression
- **Instant Messenger**
 - social networks
- **Wireless messaging**
 - semi-synchronous
- **Internet commerce**
 - the dot.com boom/bust
 - EBay, Amazon
- **Spam, phishing, ...**
 - anti-social behavior

Google™

the official **Kerry-Edwards** blog



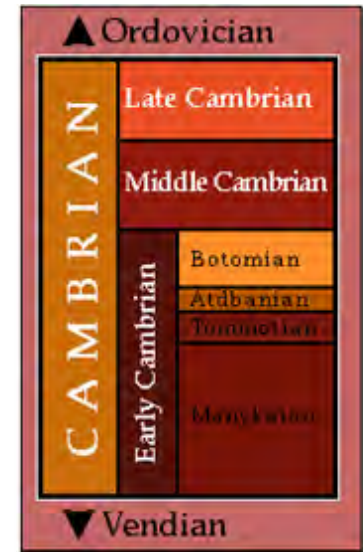
amazon.com.



ebay®

The Six Computing Eras

- **Big Iron (post WW II)**
 - vacuum tubes and campy science fiction movies
- **Mainframe ('60s/'70s)**
 - spinning tapes and bad science fiction movies
- **Workstations ('70s/'80s)**
 - spinning disks and Star Trek™
- **PCs ('80s/'90s)**
 - spinning CDs and Jurassic Park™
- **Internet ('90s)**
 - spinning DVDs and Internet pet food companies ☹
- **Implicit computing (21st century)**
 - iPods™ and The Matrix™
 - embedded intelligence in everyday objects
 - number of processors/person → infinity



Human-Computer Symbiosis



It seems reasonable to envision, for a time 10 or 15 years hence, a 'thinking center' that will incorporate the functions of present-day libraries together with anticipated advances in information storage and retrieval.

The picture readily enlarges itself into a network of such centers, connected to one another by wide-band communication lines and to individual users by leased-wire services. In such a system, the speed of the computers would be balanced, and the cost of the gigantic memories and the sophisticated programs would be divided by the number of users.

J.C.R. Licklider, 1960

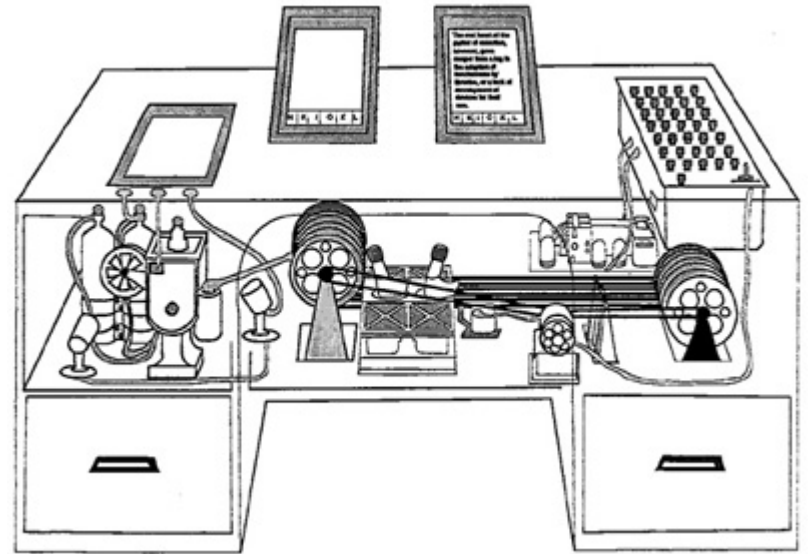


Memex: Still Prescient

“Consider a future device for individual use, which is a sort of mechanized private file and library. It needs a name, and to coin one at random, “memex” will do. A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory.”

Vannevar Bush

“As We May Think,” 1945



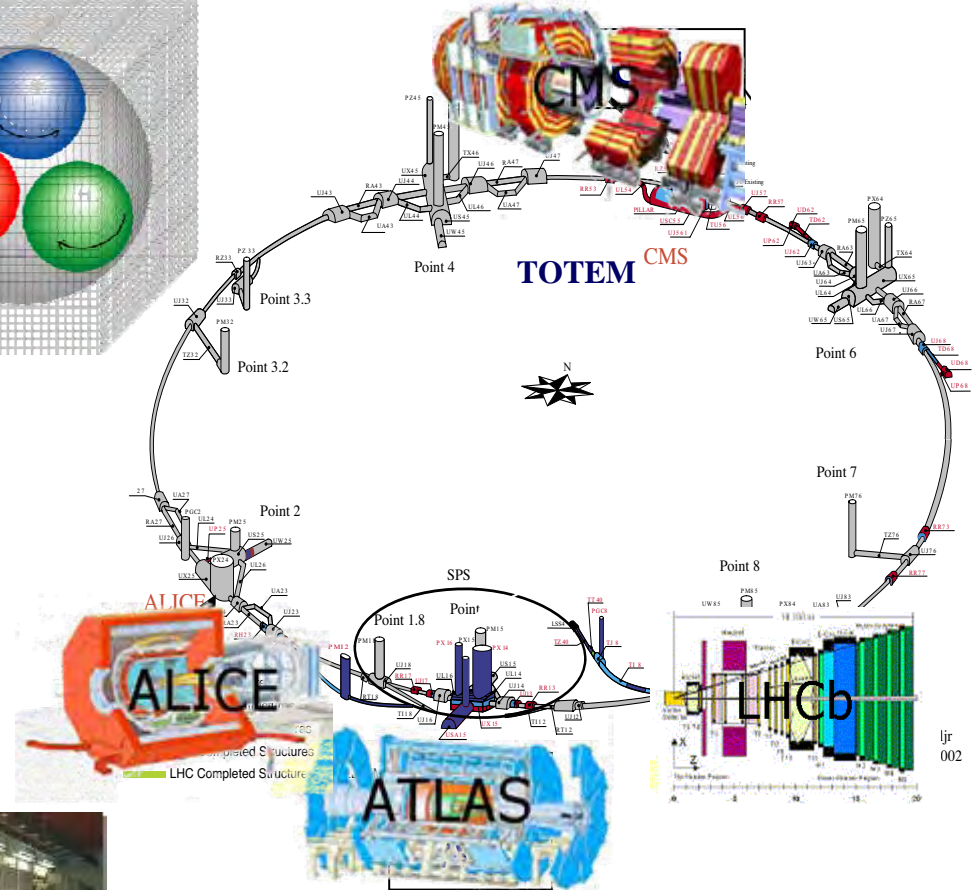
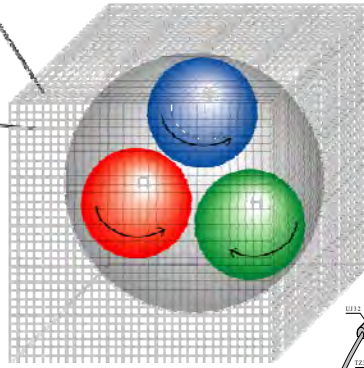
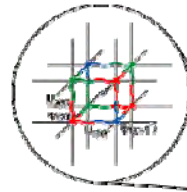
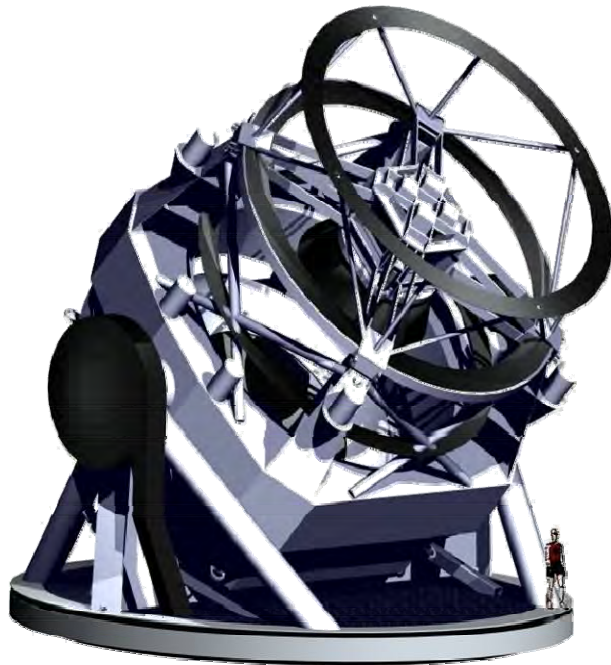
Exemplar 21st Century Challenges

- **Population growth**
 - severe weather sensitivity
 - statewide impact
 - geobiology and environment
 - economics and finance
 - sociology and policy
- **Economic research case**
 - longitudinal public health data
 - environmental interactions
 - genetic susceptibility
 - heart disease, cancer, Alzheimer's
 - privacy and insurance
 - public policy and coordination

consilience



LSST, LHC and QCD



$$D\psi(x) = \frac{1}{2a} \sum_{\mu} \left[U(x) \psi(x + \hat{\mu}) - U^{\dagger}(x - \hat{\mu}) \psi(x - \hat{\mu}) \right]$$

PITAC Report

- **Computational Science: Ensuring America's Competitiveness**
 1. *A Wake-up Call: The Challenges to U.S. Preeminence and Competitiveness*
 2. *Medieval or Modern? Research and Education Structures for the 21st Century*
 3. *Multi-decade Roadmap for Computational Science*
 4. *Sustained Infrastructure for Discovery and Competitiveness*
 5. *Research and Development Challenges*
- **Two key appendices**
 - *Examples of Computational Science at Work*
 - *Computational Science Warnings – A Message Rarely Heeded*

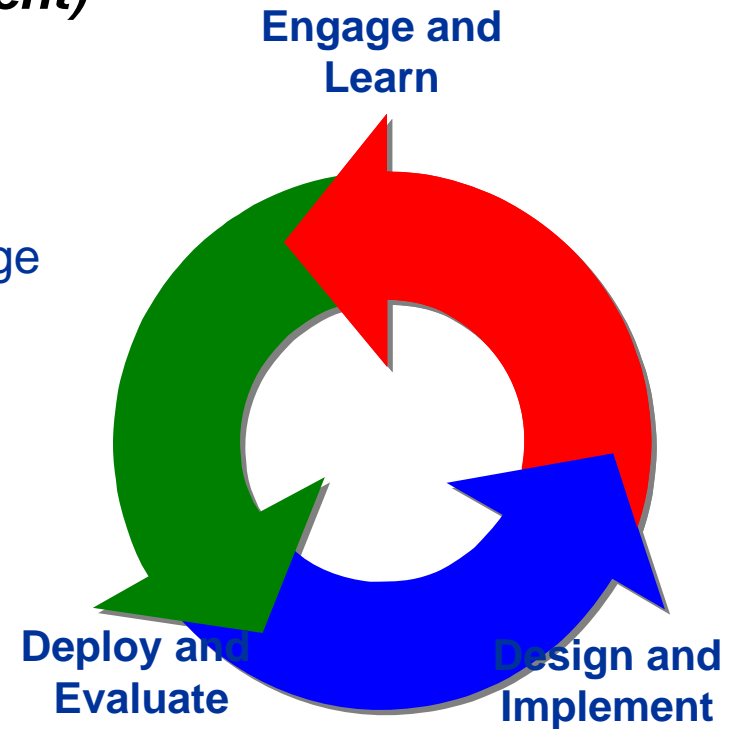


Available at www.nitrd.gov



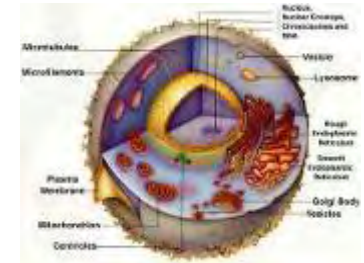
The Virtuous Cycle

- **Live in the future (*research and development*)**
 - track evolving infrastructure trends
 - prototype advanced infrastructure
- **Ride the exponentials**
 - see qualitative change from qualitative change
 - recognize the inflection points
 - e.g., personal petabytes are in sight
- **Bring insights to science (*infrastructure*)**
 - translate prototypes into use
 - glean insights from science applications
 - expand the user community
- **Learn from experience (good and bad)**
 - enhance the good
 - fix the bad and explore alternatives
- **Recognize that the cycles continue *ad infinitum***
 - commit to continual investment



Some Grand Challenges

- 1. Ubiquitous invisibility**
 - successful technologies become “invisible”
 - composable, interoperable systems
- 2. Intelligence amplification (Memex)**
 - the right information at the right time
 - seamless modality transduction, situated and mobile
- 3. Predictive *in silico* biological models**
 - the “other” artificial life
 - multidisciplinary modeling and integration
- 4. The Universe in a Box**
 - origins and alternatives
 - the theory of everything (TOE)
- 5. The Cultural Encyclopedia**
 - cultural history, context and the digital village
- 6. Grand AI, our long-term CS fascination**
 - deep questions about thinking



The Cambrian Explosion

- **Most phyla appear**
 - sponges, archaeocyathids, brachiopods
 - trilobites, primitive mollusks, echinoderms
- **Indeed, most appeared quickly!**
 - Tommotian and Atdbanian
 - as little as five million years
- **Lessons for computing and science**
 - it doesn't take long when conditions are right
 - raw materials and environment
 - leave fossil records if you want to be remembered!

