**Computing and Data Infrastructure for Large-Scale Science**

# NERSC and the DOE Science Grid:

## Support for Computational and Data-Intensive Science

*Bill Johnston*
*Distributed Systems Dept.*

*Keith Jackson*
*Distributed Systems Dept.*

*Bill Kramer*
*NERSC Center*

*Howard Walter*
*NERSC Center*

# *Outline*

- The Need for Science Grids

- What are Grids

- It Takes a Lot of Work to Make a Production Grid

- What is Missing for High-Performance Computing and Data Grids

- State of Grids

- DOE Science Grid

- Possible Roadmap for NERSC and Grids

- LBNL R&D Grid Work Couples to NERSC

- Conclusions

# ➢ *The Need for Science Grids*

- The nature of how large-scale science is done is changing
  - Distributed data, computing, people, instruments
  - Instruments integrated with large-scale computing and data systems

- "Grids" are middleware designed to facilitate the <u>routine interactions</u> of all of these resources in order to support widely distributed, multi-institutional science and engineering.
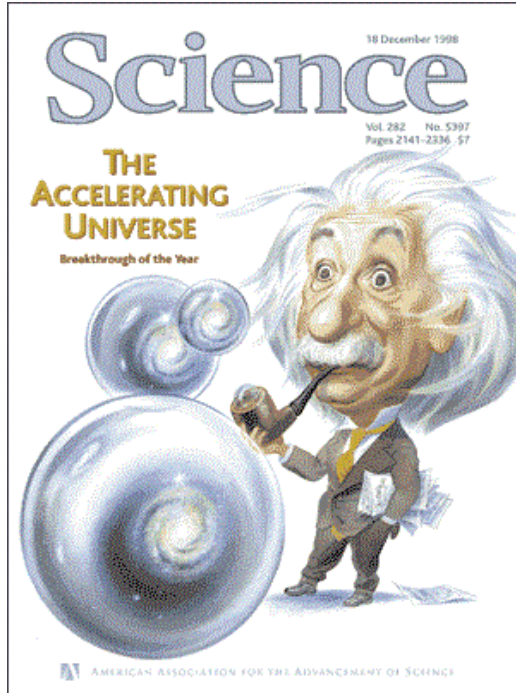
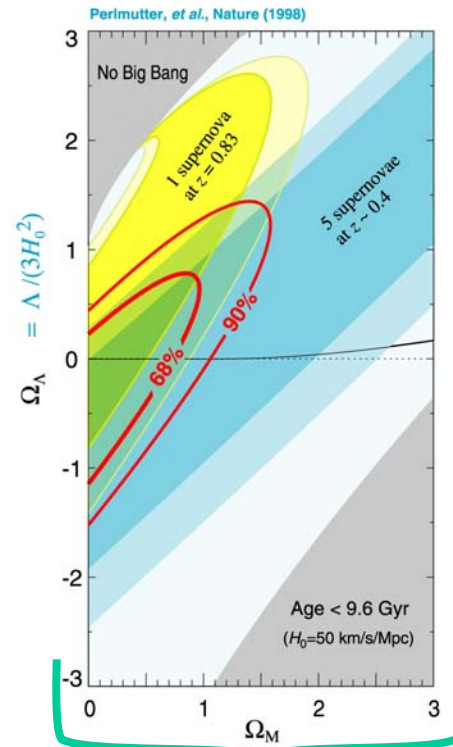# *Distributed Science Example: Supernova Cosmology*

- "Supernova cosmology" is cosmology that is based on finding and observing special types of supernova during the few weeks of their observable life

- It has led to some remarkable science, but is rapidly becoming limited by the ability of the researchers to manage the complex data-computing-instrument interactions
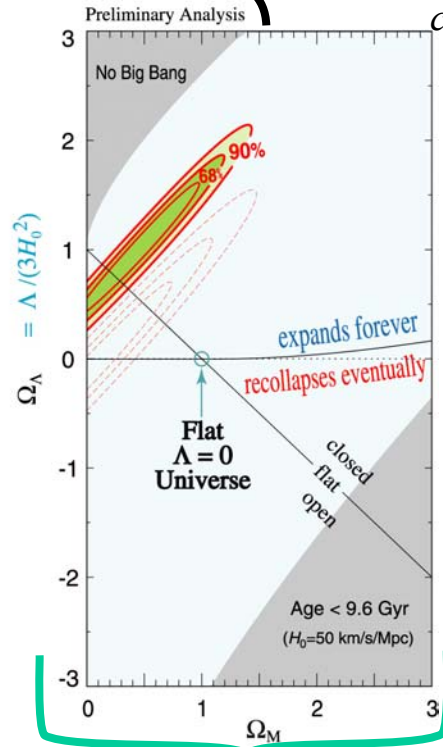
D
O
SCIENCE GRID

18 December 1998

**Science**

Vol. 282   No. 5397
Pages 2141-2336 · $7

**THE ACCELERATING UNIVERSE**

Breakthrough of the Year

AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE

Results: $\Omega$ vs $\Lambda$ from 6 supernovae

Perlmutter, et al., Nature (1998)

No Big Bang

1 supernova at z = 0.83

5 supernovae at z ~ 0.4

90%

68%

Age < 9.6 Gyr
($H_0$=50 km/s/Mpc)

$\Omega_\Lambda = \Lambda /(3H_0^2)$

$\Omega_M$

Results: $\Omega$ vs $\Lambda$ from 40 supernovae

Preliminary Analysis

No Big Bang

90%
68%

expands forever
recollapses eventually

Flat
$\Lambda = 0$
Universe

closed
flat
open

Age < 9.6 Gyr
($H_0$=50 km/s/Mpc)

$\Omega_\Lambda = \Lambda /(3H_0^2)$

$\Omega_M$
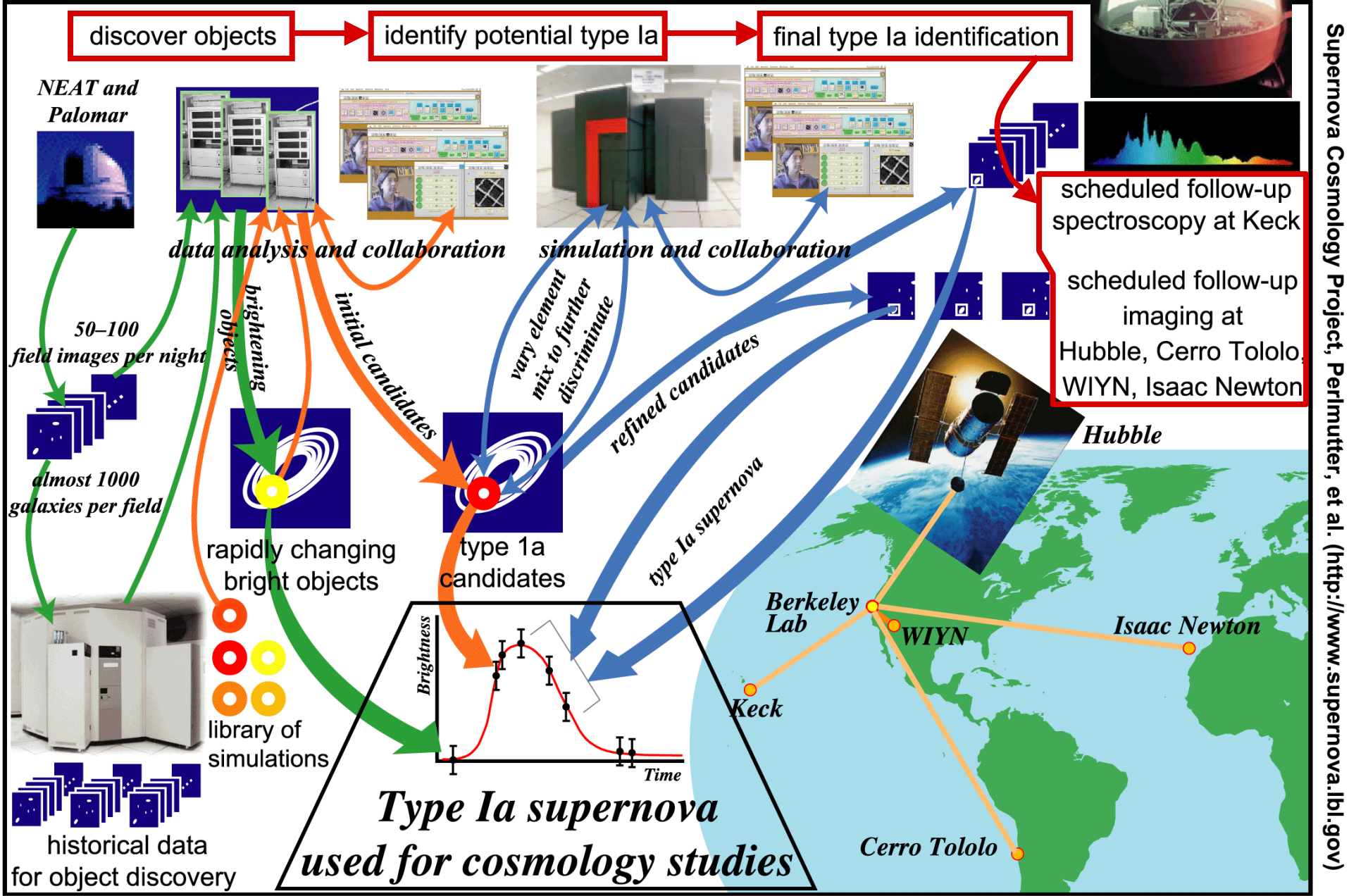
$\Lambda = \dfrac{vacuum\ energy\ density}{critical\ energy\ density}$

$\Omega = \dfrac{matter\ density}{critical\ density}$

?

| | discovered per year | | | | | | |
|---|---|---|---|---|---|---|---|
| high red shift supernova | 0 | 1 | 5 | 20 | ~20 | ~20 | ~20 |
| low red shift supernova | a few | 6 | 8 | ~10 | ~10 | ~20 | >100 |
| | *pre-1990* | *1992* | *1994* | *1996* | *1998* | *2000* | *post- 2000* |

# Supernova Cosmology Requires Complex, Widely Distributed Workflow Management



discover objects → identify potential type Ia → final type Ia identification

*NEAT and Palomar*

*data analysis and collaboration*

*simulation and collaboration*

scheduled follow-up spectroscopy at Keck

scheduled follow-up imaging at Hubble, Cerro Tololo, WIYN, Isaac Newton

*50–100 field images per night*

*almost 1000 galaxies per field*

*brightening objects*

*initial candidates*

*vary element mix to further discriminate*

*refined candidates*

*type Ia supernova*

rapidly changing bright objects

type 1a candidates

library of simulations

historical data for object discovery

*Brightness*

*Time*

## Type Ia supernova used for cosmology studies

*Hubble*

*Berkeley Lab*

*WIYN*

*Isaac Newton*

*Keck*
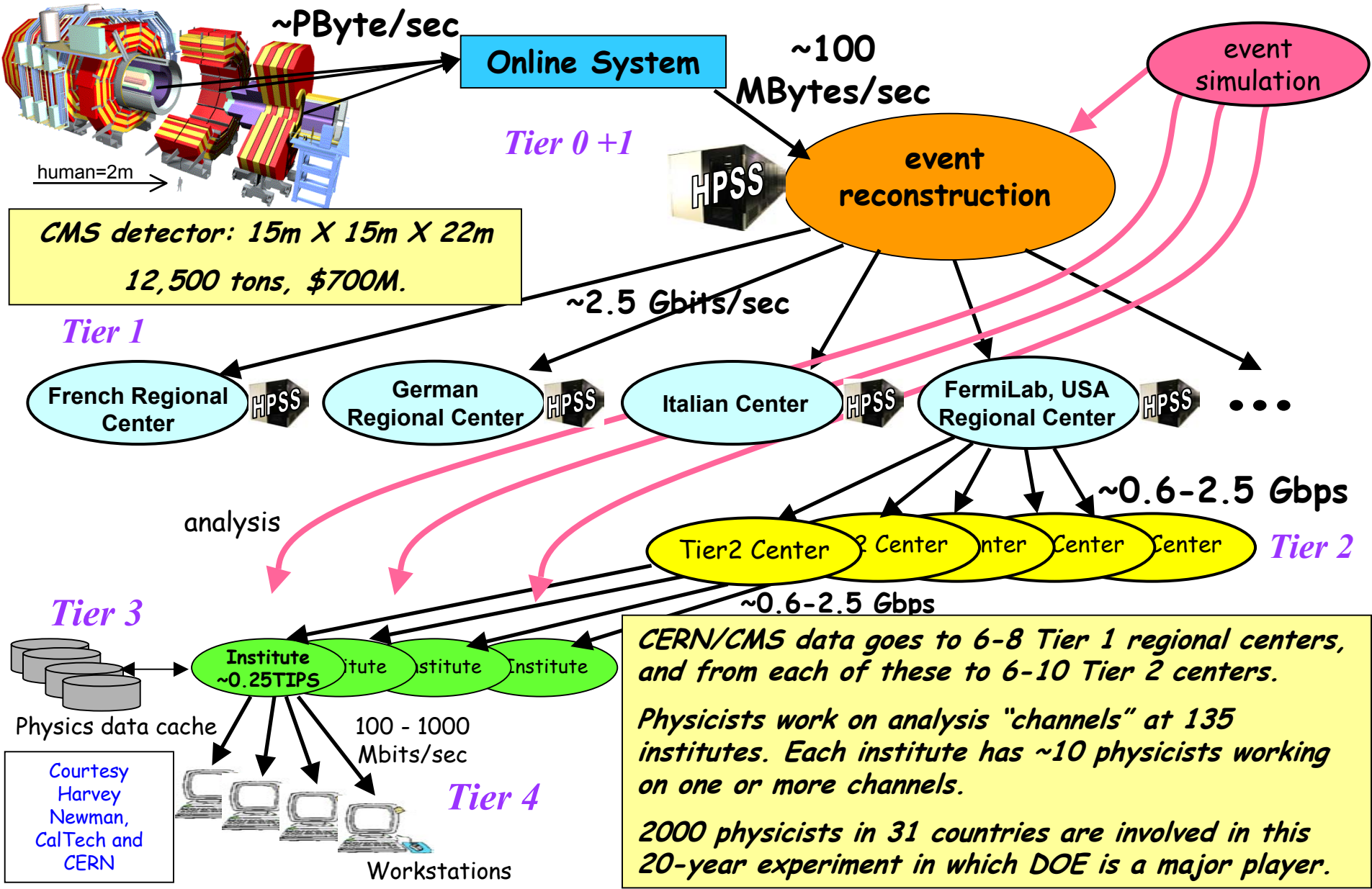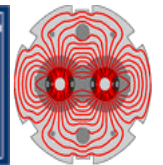
*Cerro Tololo*

# *Supernova Cosmology*

- This is one of the class of problems that Grids are focused on. It involves:

  - Management of complex workflow

  - Reliable, wide-area, high-volume data management

  - Inclusion of supercomputers in time-constrained scenarios

  - Easily accessible pools of computing resources

  - Eventual inclusion of instruments that will be semi-automatically retargeted based on data analysis and simulation

  - Next generation will generate vastly more data (from SNAP - satellite based observation)

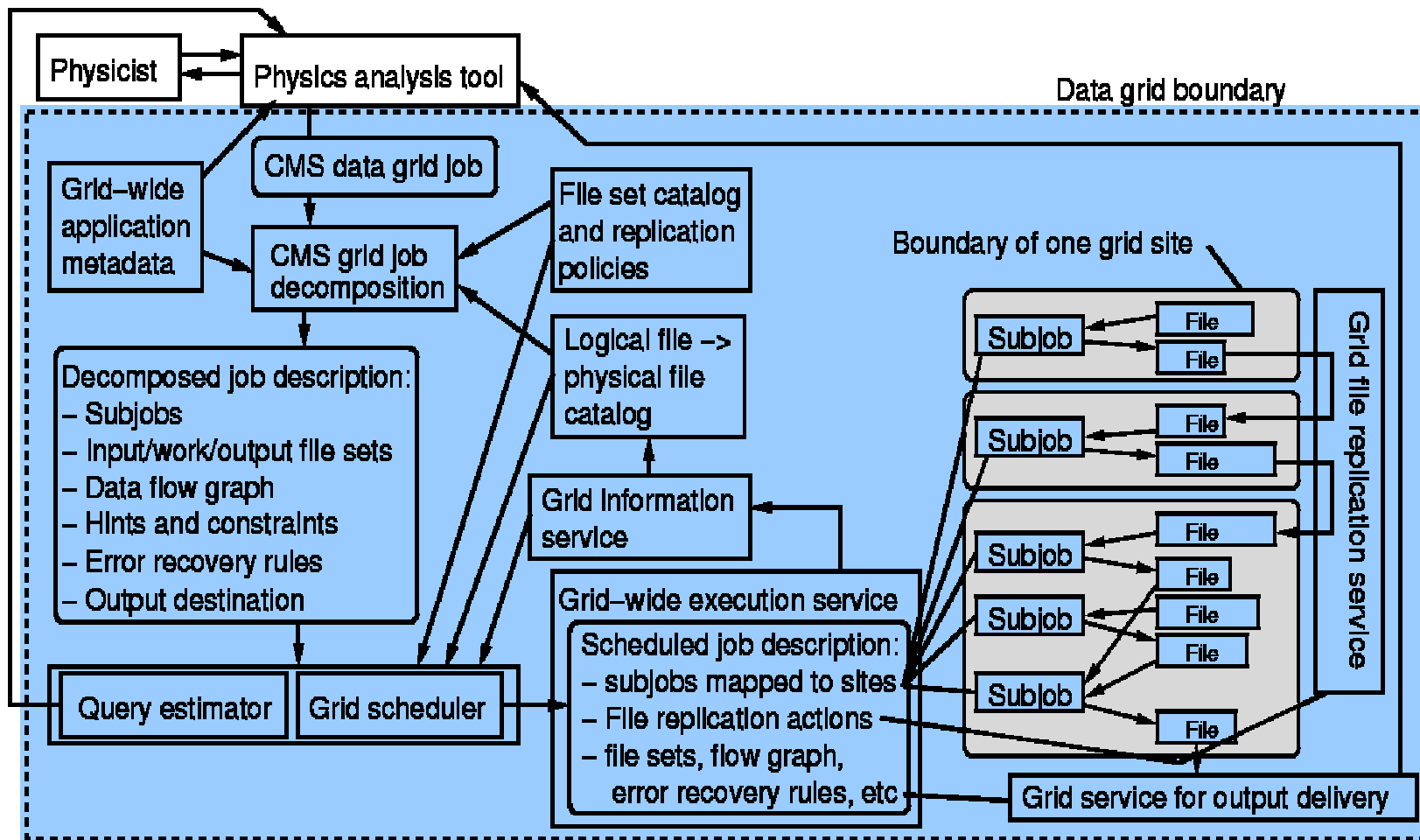# LHC Data Grid Hierarchy
## CMS as example, Atlas is similar

**~PByte/sec**

**Online System**

**~100 MBytes/sec**

**Tier 0 +1**

**event simulation**

HPSS

**event reconstruction**

CMS detector: 15m X 15m X 22m

12,500 tons, $700M.

**Tier 1**

**~2.5 Gbits/sec**

human=2m

**French Regional Center**  HPSS  **German Regional Center**  HPSS  **Italian Center**  HPSS  **FermiLab, USA Regional Center**  HPSS  **• • •**

analysis

**~0.6-2.5 Gbps**

**Tier 3**

Tier2 Center  2 Center  nter  Center  Center  **Tier 2**

**~0.6-2.5 Gbps**

Physics data cache

**Institute ~0.25TIPS**  stitute  stitute  Institute

CERN/CMS data goes to 6-8 Tier 1 regional centers, and from each of these to 6-10 Tier 2 centers.

Courtesy Harvey Newman, CalTech and CERN

100 - 1000 Mbits/sec

**Tier 4**

Physicists work on analysis "channels" at 135 institutes. Each institute has ~10 physicists working on one or more channels.

2000 physicists in 31 countries are involved in this 20-year experiment in which DOE is a major player.
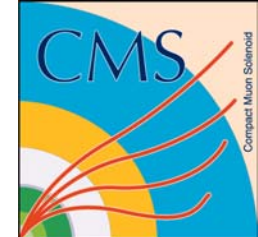
Workstations

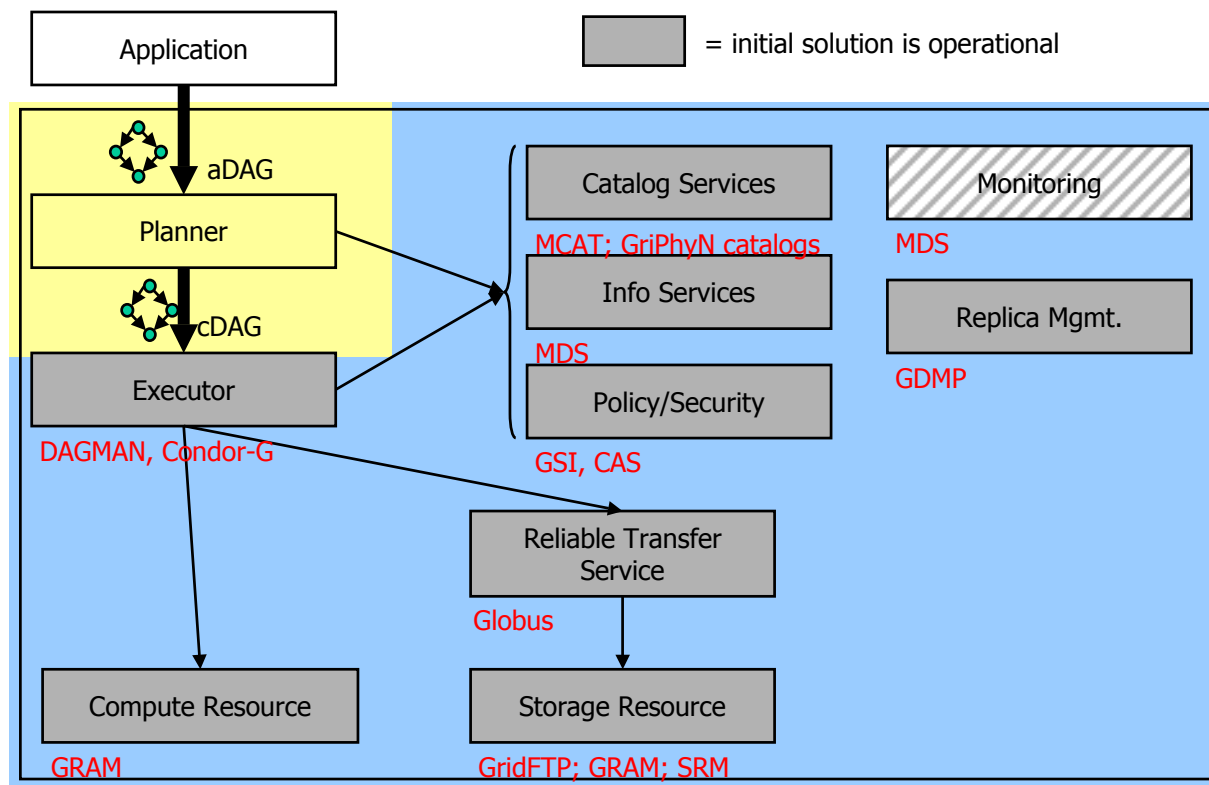# CMS Grid Requirements



Officially adopted by CMS: CMS Note 2001/037
GRIPHYN 2001-1

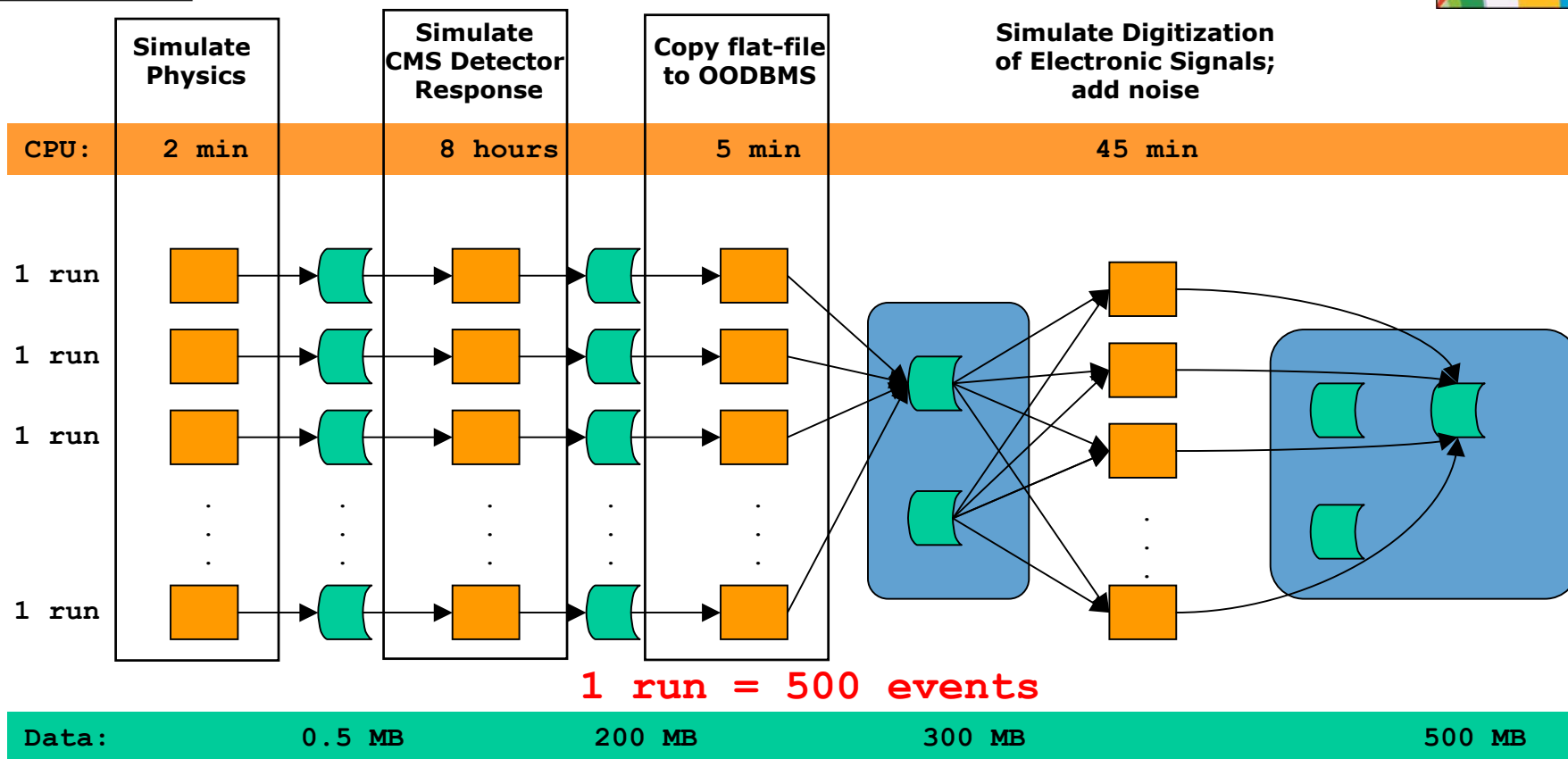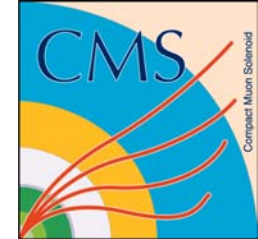# Preliminary GriPhyN
# Data Grid Architecture

- **Abstract DAGs**
  - Resource locations unspecified
  - File names are logical
  - Data destinations unspecified

- **Concrete DAGs**
  - Resource locations determined
  - Physical file names specified
  - Data delivered to and returned from physical locations

- Translation is the job of the "planner"



Application

aDAG

Planner

cDAG

Executor
DAGMAN, Condor-G

Compute Resource
GRAM

= initial solution is operational

Catalog Services
MCAT; GriPhyN catalogs

Info Services
MDS

Policy/Security
GSI, CAS

Monitoring
MDS

Replica Mgmt.
GDMP

Reliable Transfer Service
Globus

Storage Resource
GridFTP; GRAM; SRM

## Maps very well onto the CMS Requirements!

# Production of Simulated CMS Data



| CPU: | Simulate Physics | Simulate CMS Detector Response | Copy flat-file to OODBMS | Simulate Digitization of Electronic Signals; add noise |
|------|------------------|-------------------------------|--------------------------|--------------------------------------------------------|
| | 2 min | 8 hours | 5 min | 45 min |

1 run = 500 events

| Data: | 0.5 MB | 200 MB | 300 MB | 500 MB |
|-------|--------|--------|--------|--------|

- **IMPALA/BOSS (developed by CMS)**
  - **Set of scripts for mass production of simulation data**
  - **Provides parameter control and job tracking**
  - **Works quite well; produced > 20 million events!**
  - **Does not employ virtual data**

- **MOP (developed by PPDG)**
  - **Submits jobs to a Grid using DAGMan, Condor-G, Globus, and GDMP**

# Real-Time Operation of Scientific Instrument-Based Experiments

- Real-time operation of scientific instrument-based experiments is a vision that has driven much of our work in high-speed distributed systems over the past 15 years (which led directly to the Grid)

- Real-time data analysis should allow experimenters to be able to interact directly with the subject of the experiment rather than running the experiment "blind" and reconstructing after the fact what happened

# Real-Time Operation of Scientific Instrument-Based Experiments

- This can also allow insertion of human insight and computational guidance into a closed loop (servo) system so that the experiment may be driven toward a goal, or adapted in real time, depending on the sample's response to manipulation/observation

# Real-Time Operation of Scientific Instrument-Based Experiments

- One goal of Grids is to facilitate direct coupling of scientific experiments to large-scale computing systems for real-time data analysis and/or computation simulations of the subject phenomenon

- This requires persistent infrastructure and services for coordinated use of high-speed networks, distributed storage systems, and supercomputers

# *An Experiment in Quasi-real time Experiment Interaction and Collaborative Supercomputer Analysis of Micro-tomographic Data using Grid Services and Infrastructure*

Advanced Photon Source



real-time

data collection*

wide-area
dissemination

bandwidth*

real-time experiment control
and collaboration



supercomputer*
for tomographic
reconstruction

Ian Foster

*must be co-scheduled

ARGONNE ✦ CHICAGO

# ➢*What are Grids*

- Grids are technology and an emerging architecture that involves several types of middleware that sits between science portals and application, and the underlying resources (compute, data, and instrument)

- Grids are also several hundred people from the US, European, and SE Asian countries working on best practice and standards at the Global Grid Forum (www.gridforum.org)

# *Grids*

- There are several different types of users of Grid services

  - Discipline scientists

  - Problem-solving system / framework / science portal developers

  - Computational tool / application writers

  - Grid system managers

  - Grid service builders

- Each of these user communities has somewhat different requirements for Grids, and the Grid services available or under development are trying to address all of these groups

# Architecture of a Grid

## Science *Portals* and Scientific *Workflow Management* Systems

### *Web Services and Portal Toolkits*
*Applications* (Simulations, Data Analysis, etc.)
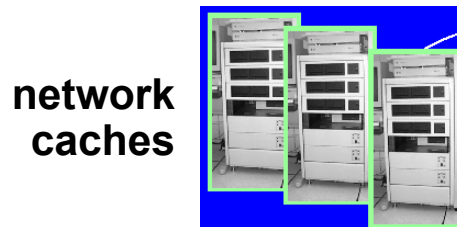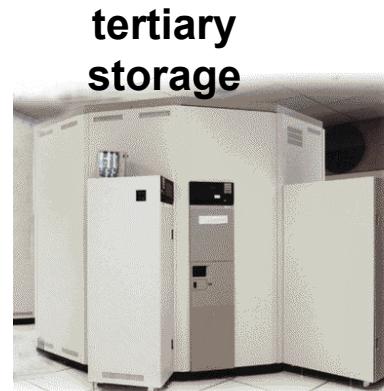*Application Toolkits* (Visualization, Data Publication/Subscription, etc.)
*Execution support and Frameworks* (Globus MPI, Condor-G, CORBA-G)

## *Grid Common Services*: Standardized Services and Resources Interfaces

| Grid Information Service | Uniform Resource Access | Uniform Data Access | Brokering | Global Event Services | Global Queuing | Co-Scheduling | Data Management | Collaboration and Remote Instrument Services | Network Cache | Communication Services | Authentication Authorization | Security Services | Auditing | Monitoring | Fault Management |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

= operational services (Globus, SRB)

## *Distributed Resources*

clusters

scientific instruments

national supercomputer

**Condor pools of workstations**

tertiary storage

network caches

## *High Speed Communication Services*

# ➤ *State of Grids*

- Grids are real, and they are useful now

- Basic Grid services are being deployed to support uniform and secure access to computing, data, and instrument systems that are distributed across organizations

- Grid execution management tools (e.g. Condor-G) are being deployed

- Data services providing uniform access to tertiary storage systems and global metadata catalogues are being deployed

- Web services supporting application frameworks and science portals are being prototyped

# *State of Grids*

- Persistent infrastructure is being built

  - Grid services are being maintained on the compute and data systems in prototype production Grids

  - Cryptographic authentication supporting single sign-on is provided through Public Key Infrastructure

  - Resource discovery services are being maintained (Grid Information Service – distributed directory service)

- This is happening, e.g., in the DOE Science Grid, EU Data Grid, UK eScience Grid, NASA's IPG, etc.

In November 2000 the Director General of Research Councils, **Dr John Taylor, announced £98M funding for a new UK e-Science programme.** The allocations were £3M to the ESRC, £7M to the NERC, £8M each to the BBSRC and the MRC, £17M to EPSRC and £26M to PPARC. In addition, £5M was awarded to CLRC to 'Grid Enable' their experimental facilities and £9M was allocated towards the purchase of a new Teraflop scale HPC system. A sum of £15M was allocated to a Core e-Science Programme, a cross-Council activity to develop and broker generic technology solutions and generic middleware to enable e-Science and form the basis for new commercial e-business software. The £15M funding from the OST for the core e-Science Programme has been enhanced by an allocation of a further £20M from the CII Directorate of the DTI which will be matched by a further £15M from industry.

What is meant by e-Science? In the future, **e-Science will refer to the large scale science that will increasingly be carried out through distributed global collaborations enabled by the Internet.** Typically, a feature of such collaborative scientific enterprises is that they will require access to very large data collections, very large scale computing resources and high performance visualisation back to the individual user scientists.

The World Wide Web gave us access to information on Web pages written in html anywhere on the Internet. **A much more powerful infrastructure is needed to support e-Science.** Besides information stored in Web pages, scientists will need easy access to expensive remote facilities, to computing resources - either as dedicated Teraflop computers or cheap collections of PCs - and to information stored in dedicated databases.

**The Grid is an architecture proposed to bring all these issues together and make a reality of such a vision for e-Science.**

# EU Data Grid

**Plans for the next generation of network-based information-handling systems took a major step forward when the European Union's Fifth Framework Information Society Technologies programme concluded negotiations to fund the Data Grid research and development project.** The project was submitted to the EU by a consortium of 21 bodies involved in a variety of sciences, from high-energy physics to Earth observation and biology, as well as computer sciences and industry. CERN is the leading and coordinating partner in the project.

**Starting from this year, the Data Grid project will receive in excess of €9.8 million for three years to develop middleware (software) to deploy applications on widely distributed computing systems.** In addition to receiving EU support, the enterprise is being substantially underwritten by funding agencies from a number of CERN's member states. Due to the large volume of data that it will produce, CERN's LHC collider will be an important component of the Data Grid

As far as CERN is concerned, this programme of work will integrate well into the computing testbed activity already planned for the LHC. Indeed, the model for the distributed computing architecture that Data Grid will implement is largely based on the results of the MONARC (Models of Networked Analysis at Regional Centres for LHC experiments) project. CERN's part in the Data Grid project will be integrated into its ongoing programme of work and will be jointly staffed by EU- and CERN-funded personnel.

# INFORMATION POWER GRID
## Engineering and Research Site

NASA   IPG

| IPG News | Support | Engineering | Research | About IPG | Launch Pad |

| Concepts | Vision | Presentations | Workshops | Reviews |

## What is the Information Power Grid?

The **Information Power Grid** (IPG) is NASA's high performance computational grid. Computational grids are are persistent networked environments that integrate geographically distributed supercomputers, large databases, and high end instruments. These resources are managed by diverse organizations in widespread locations, and shared by researchers from many different institutions.

The IPG is a collaborative effort between NASA Ames, NASA Glenn, and NASA Langley Research Centers, and the NSF PACI programs at SDSC and NCSA, and is funded by the Computing, Information and Communications Technology (CICT) program at NASA Ames Research Center.

http://www.ipg.nasa.gov/

# It Takes a Lot of Work to Make a Production Grid

Based on three years experience building IPG at NASA, we can enumerate ***Steps for Building a Multi-Site, Globus-Based, Computational and Data Grid:***

- To Start - build a testbed for training:
  - Establish an Engineering Working Group that involves the Grid deployment teams at each site
  - Identify the computing and storage resources to be incorporated into the Grid
  - Involve the systems and operations staff (computation and storage) as early as possible
  - Plan for a Grid Information Service (directory) server at each site with significant resources
  - Build Globus on test systems

- Preparing for the transition to a prototype-production Grid means addressing a number of significant scaling issues:

    – Set up, or identify, a Certification Authority to issue Grid PKI identity certificates to Grid users and hosts

    – Make sure that you understand the extent to which need a homogeneous Certificate Policy

    – Issue host certificates for the resources and establish procedures for installing them

    – Count on revoking and re-issuing all of the certificates at least once before going operational

    – Using certificates issued by your CA, validate correct operation of the Globus security libraries, GSI ssh, and GSI ftp

- Determine the "boundaries" of your Grid:
  - What CAs you trust establishes the maximum extent of your user population
  - The scope for searching of GIS/GIISs establishes default resource domain

- Establish the conventions for the Globus mapfile (the Grid resource authorization mechanism)

- Validate network connectivity between the sites and establish agreements on firewall issues

- Provide the application developers with end-to-end monitoring

- Build and test your Grid incrementally

- Establish user help mechanisms

- At this point, Globus and the basic infrastructure should be operational on the testbed system(s). The Globus deployment team should be familiar with the install and operation issues, and the sys admins of the target resources should be engaged.

- Next step is to build a prototype-production environment:

  – Deploy and build Globus on at least two production computing platforms at two different sites

  – Establish the relationship between Globus job submission and the local batch schedulers (one queue, several queues, a Globus queue, etc.)

- Validate the services for moving data between all of the systems involved in your Grid

- Decide on a Grid job tracking and monitoring strategy

- Put up one of the various Web portals for Grid resource monitoring

➢ Establish a Grid/Globus application specialist group

➢ Identify early users and have the Grid/Globus application specialists assist them in getting jobs running on the Grid

- ## Knowledge Frameworks

  – From a problem description formulated by a scientist or engineer, be able to identify and automatically invoke appropriate operations on the computational components and datasets of the discipline area to "solve" the problem

- ## Science Portals/Problem-Solving Environments

  – General mechanisms and toolkits are needed for representing and manipulating the structure of the problem, and for easily building portals to instantiate this (e.g. with Web services)

- ## Workflow Management

  – Provide for description and subsequent control of the related steps and events that represent a "job." A general approach is needed to provide a rule-based execution management system driven from published/subscribed global events (where the "events" represent process completion, file or other state creation, instrument turn-on, etc.)

# *What is Missing for High-Performance Computing and Data Grids*

- Collaboration frameworks
  - Mechanisms for human control and sharing of all aspects of an executing workflow

- Global File System
  - Should provide Unix file semantics, be distributed, high performance, and use the Grid Security Infrastructure for authentication

- Application "wrapping" and composing
  - Must enable dynamic object management in an environment of widely distributed resources (e.g. NSF GRADS)

- Monitoring
  - Needed for all aspects of a running job (e.g. to support fault detection and recovery)

- ## Authorization
  - Mechanisms to accommodate policy involving multiple stakeholders providing use-conditions on resources and user attributes in order to satisfy those use-conditions

- ## Dynamic construction of execution environments supporting complex distributed applications
  - Co-scheduling many resources to support transient science and engineering experiments that require combinations of instruments, compute systems, data archives, and network bandwidth at multiple locations (requires support by resource)

- ## Grid interfaces to existing frameworks: CORBA, MPI, Condor-G, DCOM

# Combined Grid and Web Services Architecture

| Clients | Application Portals | Web Services | Grid Services: Collective and Resource Access | Resources |
|---|---|---|---|---|

**Clients**

- X Windows
- Web Browser
- PDA

http, https. etc.

**Application Portals**

- Discipline / Application Specific Portals (e.g. SDSC TeleScience)
- Problem Solving Environments (AVS, SciRun, Cactus)
- Environment Management (LaunchPad, HotPage)
- composition frameworks (e.g. XCAT)

Python, Java, etc. JSPs

*Apache SOAP, .NET, etc.*

XML / SOAP over Grid Security Infrastructure

**Web Services**

- Job Submission / Control
- File Transfer
- Data Management
- Monitoring
- Events
- ……
- Credential Management
- Workflow Management
- other services:
  - visualization
  - interface builders
  - collaboration tools
  - numerical grid generators
  - etc.

CoG Kits implementing Web Services in servelets, servers, etc.

*Apache Tomcat&WebSphere &Cold Fusion=JVM + servlet instantiation + routing*

Grid Protocols and Grid Security Infrastructure

**Grid Services: Collective and Resource Access**

- Grid ssh
- CORBA
- GRAM
- Condor-G
- SRB/ Metadata Catalogue
- Data Replica and Metadata Catalog
- GridFTP
- Grid Monitoring Architecture
- Grid X.509 Certification Authority
- MPI
- Grid Information Service
- Secure, Reliable Group Comm.
- Grid Web Service Description (WSDL) & Discovery (UDDI)

Grid Protocols and Grid Security Infrastructure

**Resources**

- Compute (many)
- Storage (many)
- Communi-cation
- Instruments (various)

# *DOE Science Grid*

- SciDAC project to explore the issues for providing persistent operational Grid support in the DOE environment: LBNL, NERSC, PNNL, ANL, and ORNL

  - Initial computing resources
    - $\approx$ 10 small, medium, and large clusters
  - High-bandwidth connectivity end to end (high-speed links from site systems to ESnet gateways)
  - Storage resources: four tertiary storage systems (NERSC, PNNL, ANL, and ORNL)
  - Globus providing the Grid Common Services
  - **Collaboration with ESnet for security and directory services**

**User Interfaces**
**Application Frameworks**
**Applications**

**Grid Services:** Uniform access to distributed resources

Grid Information Service | Uniform Resource Access | Brokering | Global Queuing | Co-Scheduling | Global Event Services | Data Cataloguing | Uniform Data Access | Collaboration and Remote Instrument Services | Network Cache | Communication Services | Authentication Authorization | Security Services | Auditing | Monitoring | Fault Management

**Grid Managed Resources**

Asia-Pacific

Europe

PNNL

ESNet

MDS CA

ANL

ESnet

**DOE Science Grid**

LBNL

NERSC Supercomputing & Large-Scale Storage

ORNL

Initial Science Grid Configuration

Funded by the U.S. Dept. of Energy, Office of Science,
Office of Advanced Scientific Computing Research,
Mathematical, Information, and Computational Sciences Division

- This is a roadmap that introduces Grids into NERSC and tries to fit into the DOE constraints. We can still have a big impact on DOE science.

- Initial strategy is to use Grid services to provide services that NERSC must provide in some form in any event — this should provide a natural, gradual, and significant integration of Grids into NERSC.

- The ordering is based on considering criticality for Grids, data orientation, complexity of introduction, and maturity of the service.

## 1) Grid security and authorization infrastructure

- Single sign-on
  - User authentication via PKI identity certificate, followed by derived proxy certificates for all subsequent authorization

- Grid authorization for NERSC resources
  - Per-resource access control list for Grid users

- Establish policy agreements on characteristics of PKI identity certificates
  - This will permit interoperation with other DOE-relevant Grids

## 2) GridFTP access to NERSC production HPSS

– GridFTP is not your ordinary FTP - provides enhancements needed for high speed and high volume

- Can adapt to network
- Parallel streams for high data rates
- Reliable transfer for high data volume
- HPSS queue management and caching combined with the HRM (Hierarchical Resource Manager - HPSS manager)

- Tasks 1 and 2 are an important first step toward Grids: They provide high-speed WAN access to a Grid-enabled NERSC HPSS, thereby allowing it to be integrated with Grid applications - e.g. the SciDAC Earth Systems Grid.

## 3) Integrate Alvarez cluster into the Grid

- – DOE Science Grid funded (in part)

- – Would add Grid Information Service and the Globus job manager to NERSC expertise

- – Would provide a NERSC computational resource

  on the Grid

## 4) Pursue testing supplied Grid services on the IBM SPs

- – Minimizes the cost to NERSC for getting SPs
  onto the Grid

- – An important first step to getting the flagship NERSC systems on the Grid

**Steps 1-4 can probably be done in the NERSC Base Program.**

Further services probably require additional funding.

5) NERSC Web Grid services server:

- Makes NERSC Grid services available as Web services
- Simplifies construction of application portals/frameworks like the PNNL ECCE, Cactus, SDSC Tele-microscopy portal, etc.
- Leverages a huge commercial effort to build the tools for constructing Web portals from Web services (MS.NET, Allaire Cold Fusion, IBM WebSphere, Apache Tomcat, etc.)

In my opinion, this is one of the big payoffs for a NERSC Grid: Significantly lowering the barrier to developing complex science application frameworks that include NERSC services and resources

## 6) Quality of Service on the supercomputers

– Initially as advanced reservation of job run slots in the batch queuing system

– This provides the potential to integrate the NERSC supercomputers into time-constrained, or even real-time, experiment environments

## 7) Data Grid services (from NSF GriPhyN and EU Data Grid projects)

– Data and metadata catalogue

– Data location management

– Support for "virtual" (un-materialized) data

• Procedures in metadata catalogue

• Data generation workflow system

# ➢ *LBNL R&D Grid Work Couples to NERSC*

There is a strong Grids program in LBNL computer science. This is coupled to NERSC through the DOE Science Grid and SciDAC. For example, we have projects funded in:

- DOE Science Grid
- Self-configuring network monitoring
- Secure, reliable group communication
- Distributed security
- Grid Monitoring Architecture
- High-speed data management
- Network-aware Operating Systems
- Hierarchical storage resource manager
- Distributed collaboration services
- Grid Web Services
- End-to-end performance monitoring and analysis
- Science workflow management

Many of these are closely related to Grid Forum activities

# *LBNL R&D Grid Work Couples to NERSC*

There is a very important Grids program in ESnet

- Operate an identity Certification Authority for the science community of DOE and its collaborators
  - Greatly eases participation in Grids by science organizations

- Establish common identity Certification Authority policy
  - A major issue for the collaboration of the DOE Labs, the non-DOE HEP Labs, and Universities - ESnet is taking the lead in helping to homogenize policy in US and European science labs

- R&D in large-scale Directory Services
  - Addresses a major issue for the scalability of Grids - how do you do resource discovery as you get hundreds of organizations participating in a common Grid?

These are also involved in Grid Forum Working Groups

# Conclusions

- *NERSC can join the Grids environment in an adiabatic way initially, followed by a slow ramp-up, and still provide a significant Grids environment to its users*

- Grids have the potential to make the construction and operation of large-scale distributed science computing and data-management environments much easier than today — *this should increase the scope and productivity of the large-scale science applications*

- Grids will help *prevent continual reinvention* of distributed services