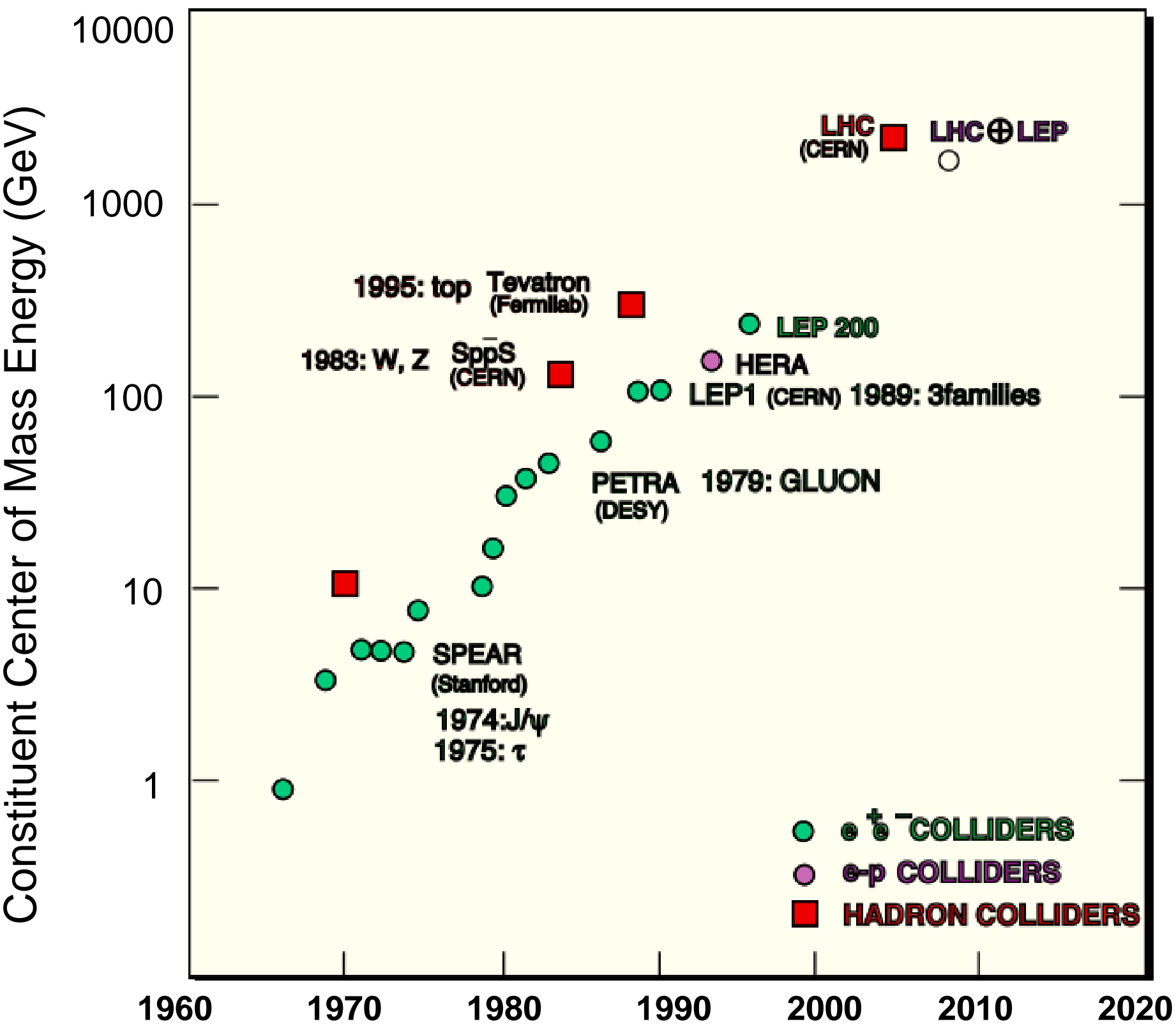
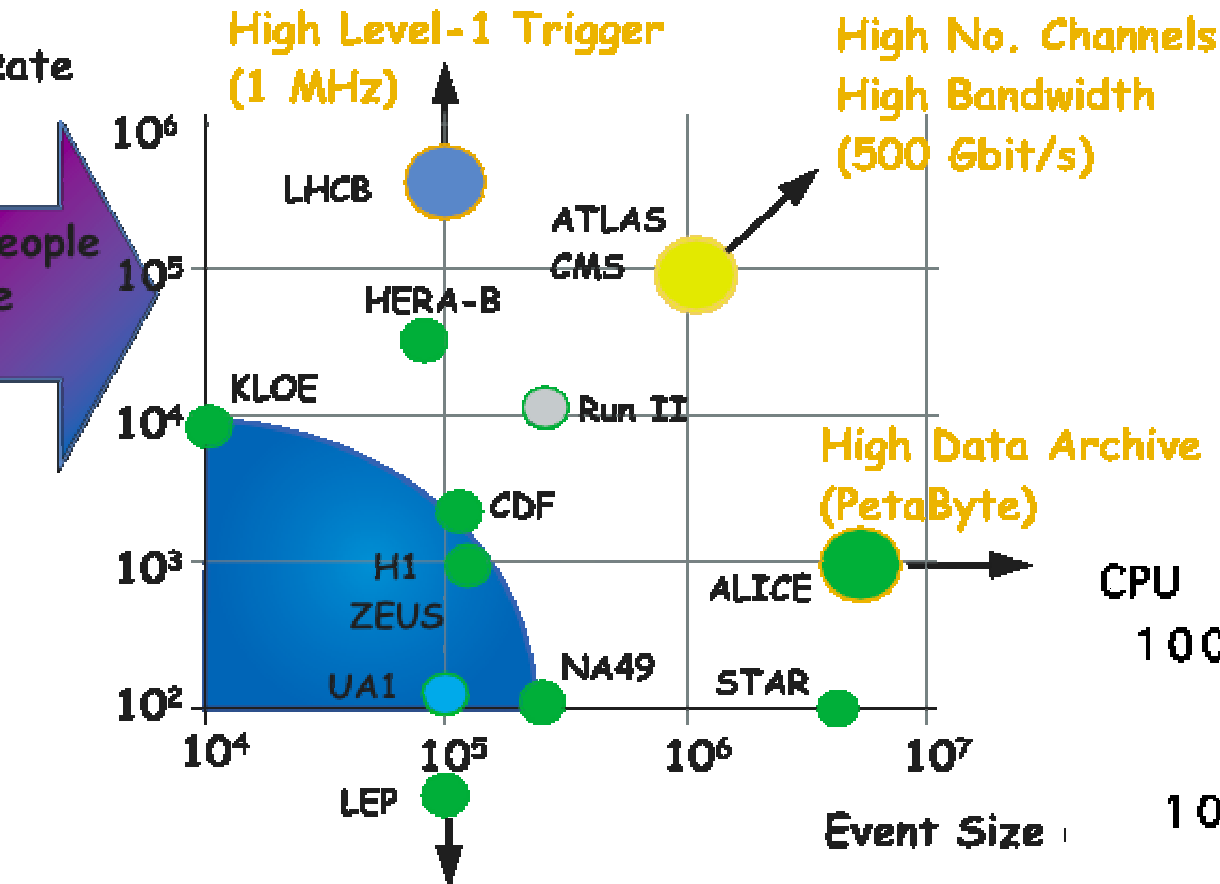


# Federating Grid Resources

Michael Ernst  
DESY Seminar  
May 3rd 2004



# In Data Rate, Data Size, CPU and Number of Scientists



CPU

100,000

10,000

1,000

100

10

0

500

1000

1500

2000

Earth Simulator

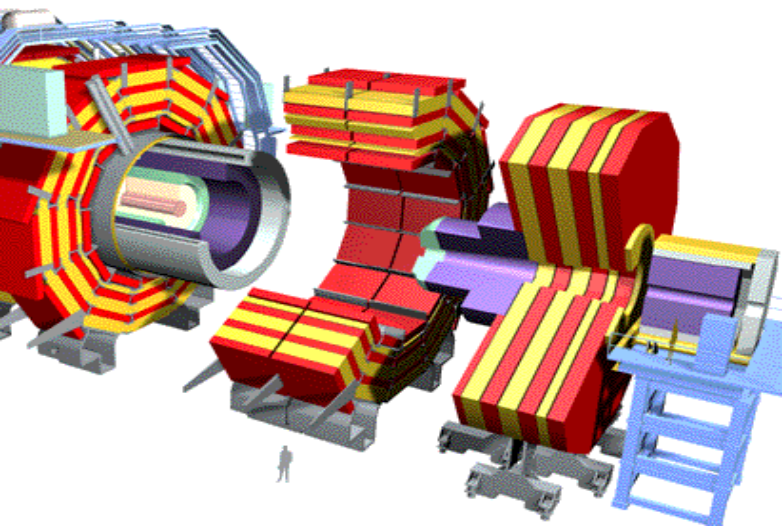
Grav. Wave

Astronomy

Atmospheric Chemistry Group

Nuclear Exp.

Current accelerator Exp.

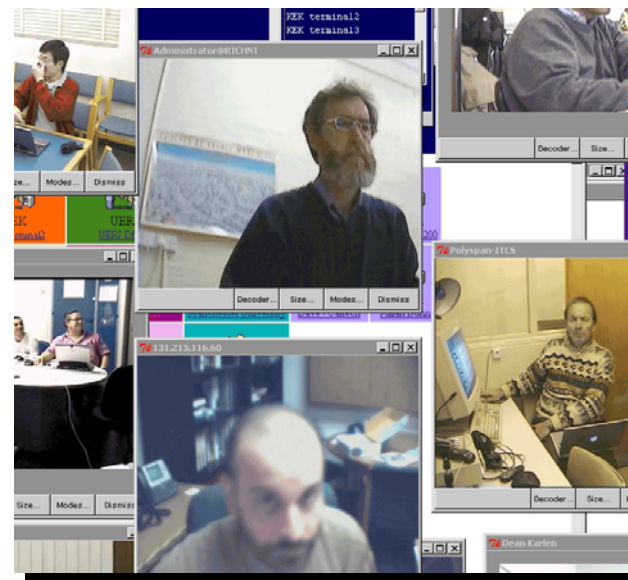


+



e HEP Computing Unprecedented in Scale and Complexity (and Costs)  
an Advanced Coherent Global “Information-Infrastructure”  
ational and Interdisciplinary Partnerships

# Over the Universities to do Research on Physics Data



why we are interested in Grids and enabling Information Technology

global collaboration of thousands of physicists

Provide capabilities to individual physicists and communities of scientists that allow

- To participate as an equal in the research program
- To be fully represented in the Global Experiment Enterprise
- To on-demand receive whatever resources and information they need to explore their science interests respecting the collaboration wide priorities and needs

massive computing, storage, networking resources

including “opportunistic” use of resources that are not owned by a particular experiment!

full access to dauntingly complex “meta-data”

That need to be kept consistent to make sense of the event data

# Architectural) Distributed Computing Model with multiple Tiers National Centers: Managed, fair-shared access to data for Physics where

maximize total funding resources while meeting the

local computing and data handling needs

balance between proximity of datasets to appropriate resources,

and to the users => Tier-N Model

efficient use of network: higher throughput

Per Flow: Local > regional > national > international

local intellectual resources, in several time zones

Laboratories, universities, remote sites

Involving physicists and students at their home institutions

greater flexibility to pursue different physics interests, priorities, and

resource allocation strategies by region

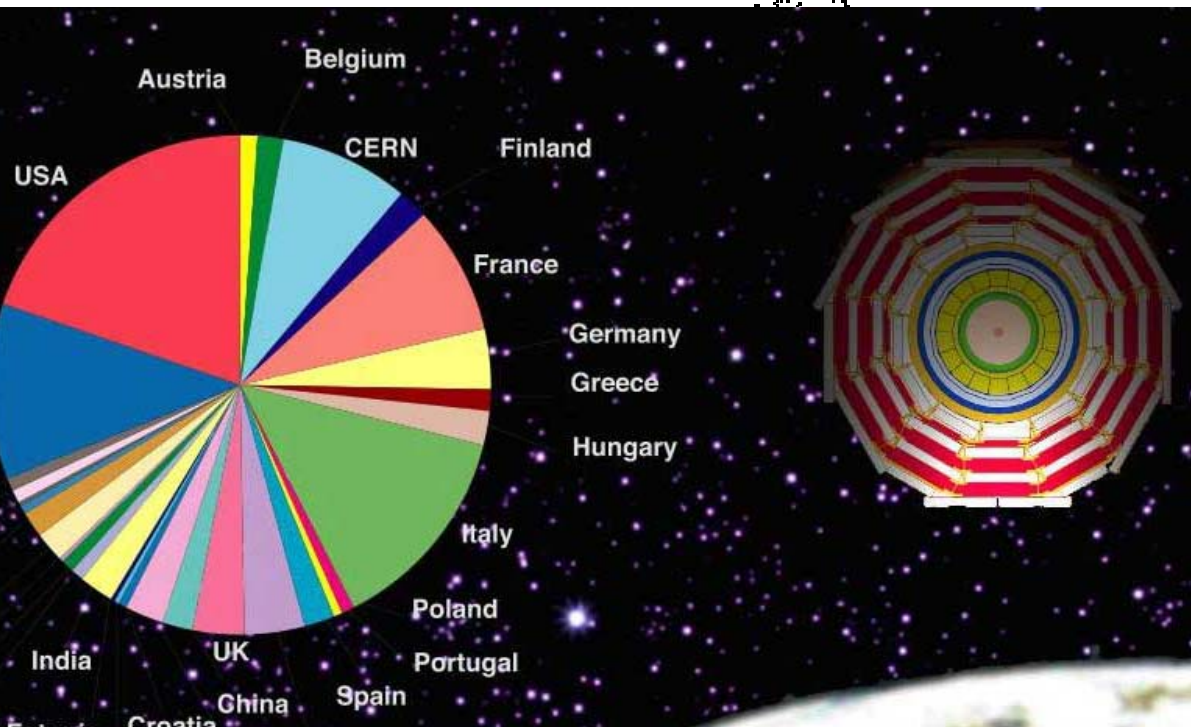
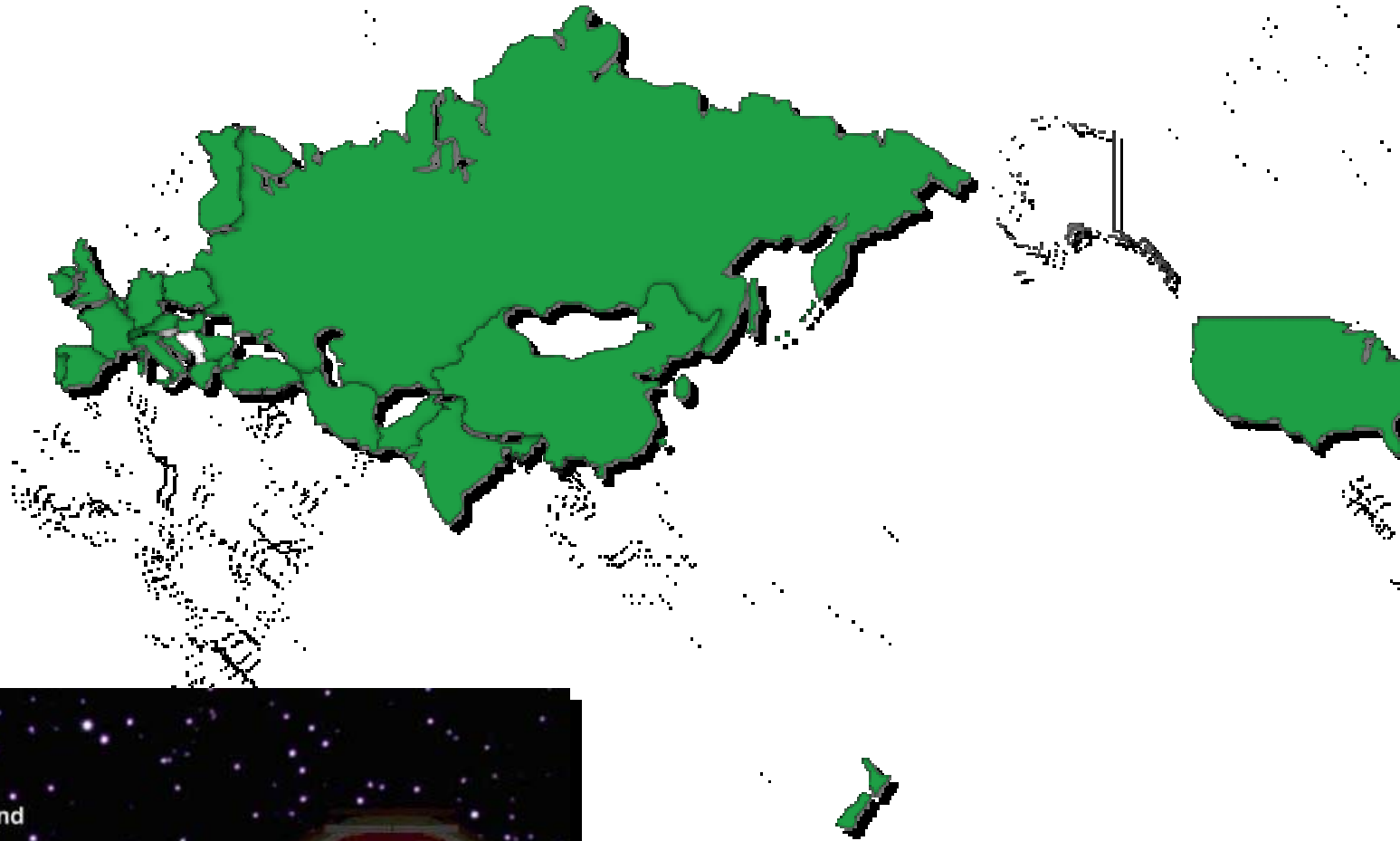
and/or by common interests (physics topics, subdetectors,...)

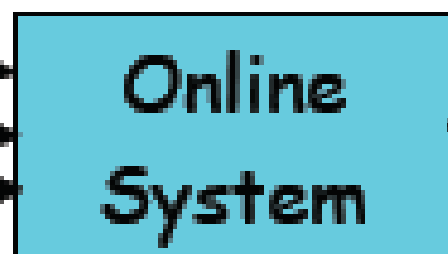
manage the System's Complexity

Partitioning facility tasks, to manage and focus resources

36 Nations, 159 Institutions, 1940 Scientists and Engineers (February 1994)

# The CMS Collaboration





100-200 MBytes/s

**Tier 0**



2.5 - 10 Gbits/s



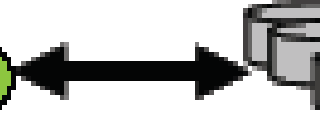
**Tier 2**



2.5

~0.6 Gbits/s

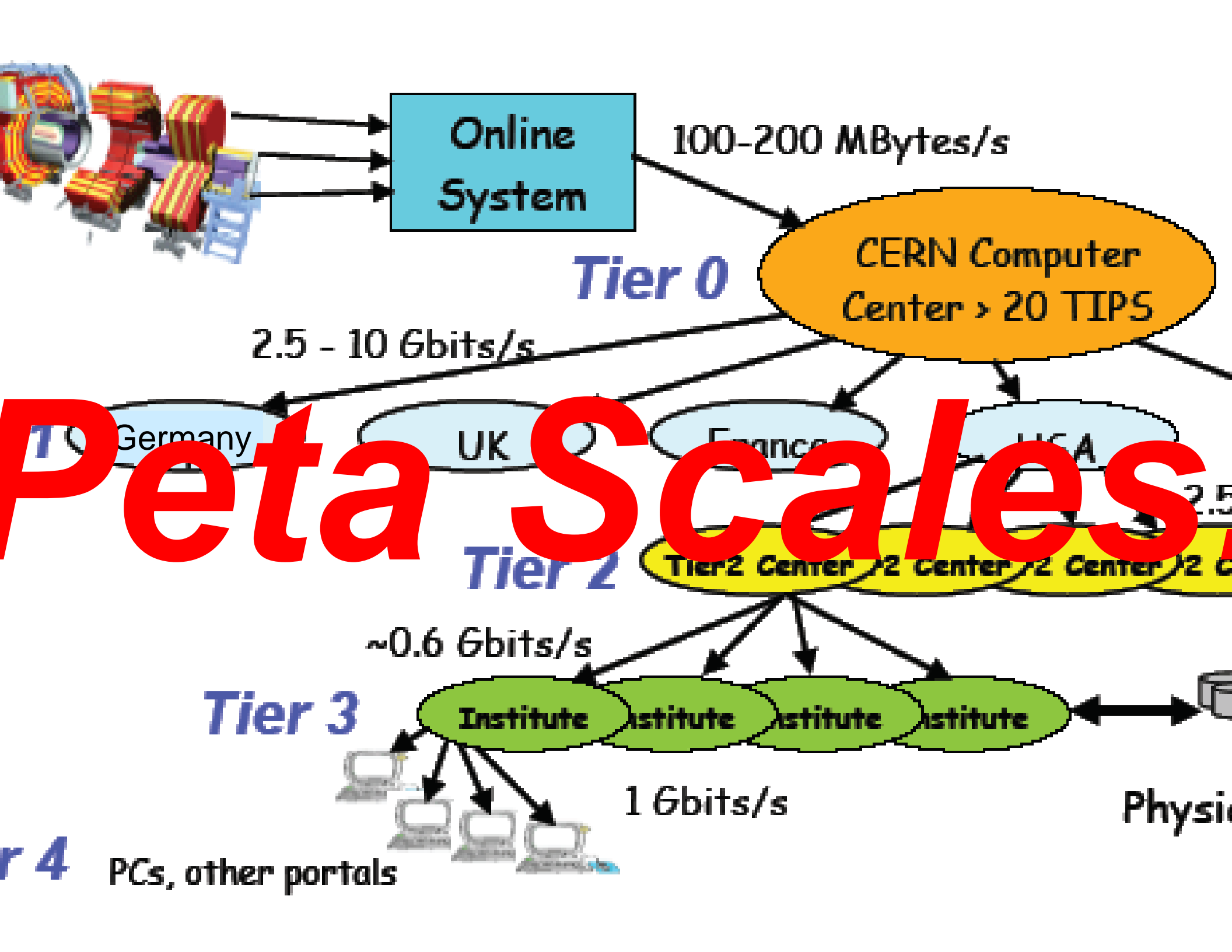
**Tier 3**



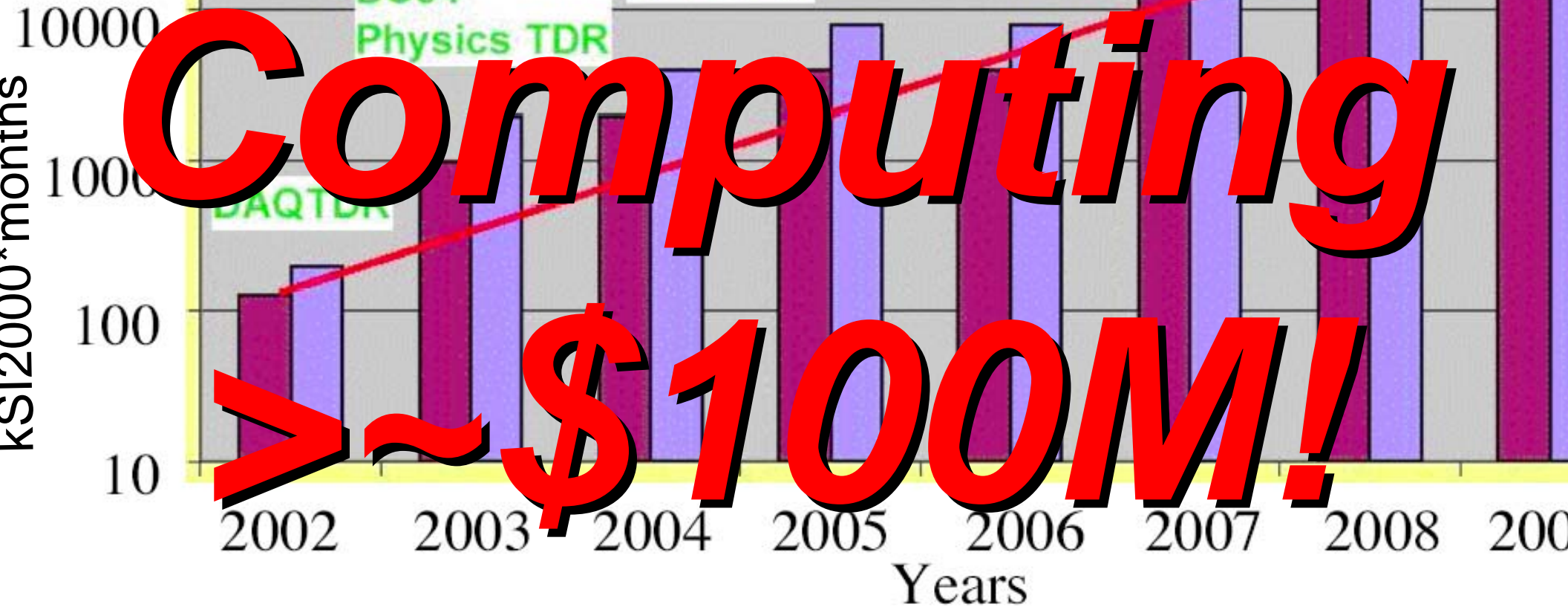
1 Gbits/s

Physic





**Total Costs For**



to enable National and Regional Organizations to meet the  
requirements

National Funding Agencies pay for computing contributions of projects  
An independent set of interfaces leads to partitioning of resources

More efficient to share large computing facilities with interoperable  
interfaces and Environments

By itself there is no intrinsic value in globally distributing computing  
resources

Success of Grid Computing is based on three fairly simple principles  
National Funding Agencies prefer to spend money in home countries  
Existing local resources (physical infrastructure, HR) can be leveraged  
plus matching funds and university grants

Computing Clusters are specified for peak needs and the usage  
has structure => Spare Computing Cycles available somewhere

Common Interfaces are prerequisite to discover available resources  
Different Communities likely to benefit from Grid Computing

Project to build a common grid environment to:

Provide the infrastructure and services needed for production and analysis applications running at scale in a common grid environment.

Provide the next phase of the International Virtual Data Grid Laboratory (iVDGL).

Provide a platform for computer science technology demonstrators.

Project between U.S. ATLAS and U.S. CMS to use a common environment at the LHC Tier-1 and Tier-2 centers:

Allowing software developed by one experiment to be integrated and used by both.  
Provides agreement on policies, principles and procedures for Grid system use.  
Enabling opportunistic use of additional non-HEP computing resources.

Goal: Demonstrate and Operate a Functioning Multi-Organization Grid:

With well-defined metrics -- a thousand running processes, TBytes/day data transfer.  
Supporting CMS, ATLAS, SDSS, LIGO, Biology, CS applications.

U.S. CMS Grid2003 is a continuation and extension of the existing U.S. CMS Grid.  
Wanted to participate in a multi-experiment and organization Grid environment.

***Development and Integration Test Grids are Essential!***



is a Collaborative Team Effort of Application Integrators and Deployers, Scientists and Grid Service Providers and Supporters. Coordinated by a Taskforce bringing the Stakeholders and joint coordinators representing iVDGL & PPDG and CMS.

Participants:

ATLAS & CMS

um Grid Projects:

International Virtual Data Grid Laboratory (iVDGL), which includes LIGO, SDSS

article Physics Data Grid Collaboratory Pilot (PPDG)

Grid Physics Network (GriPhyN)

Telemetry (University of Chicago)

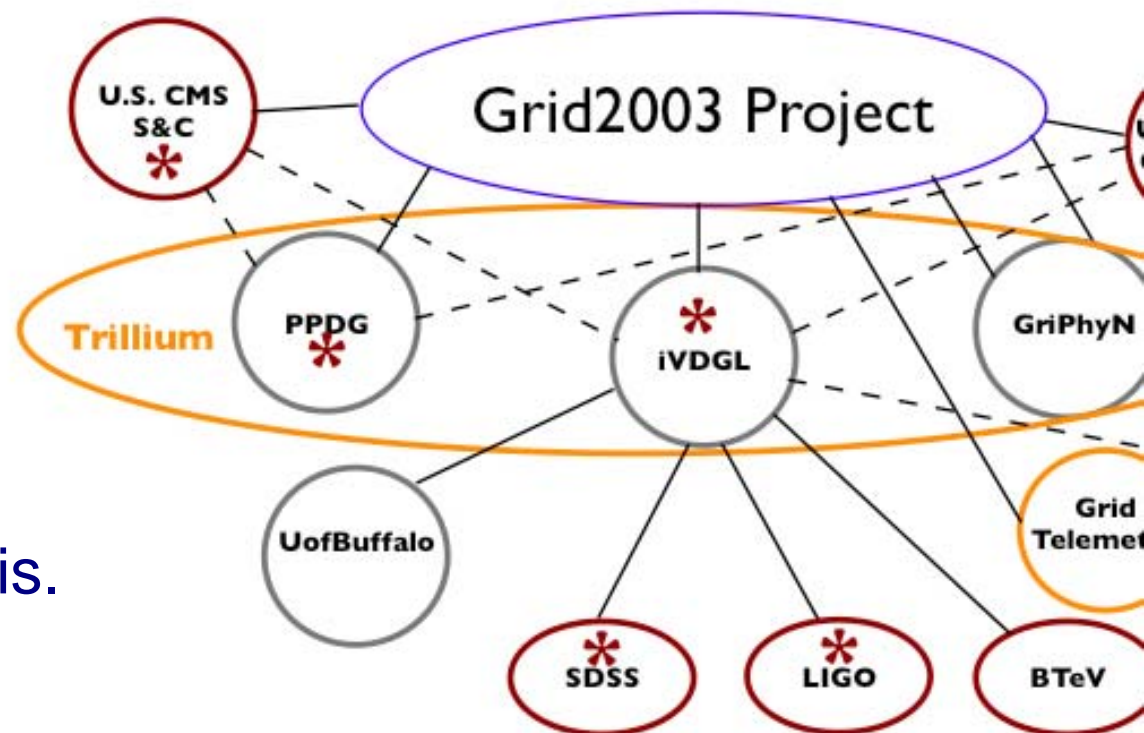
oined During the Project:

Korea

/

Biology - protein sequence analysis.

ersity of Buffalo - Biology -



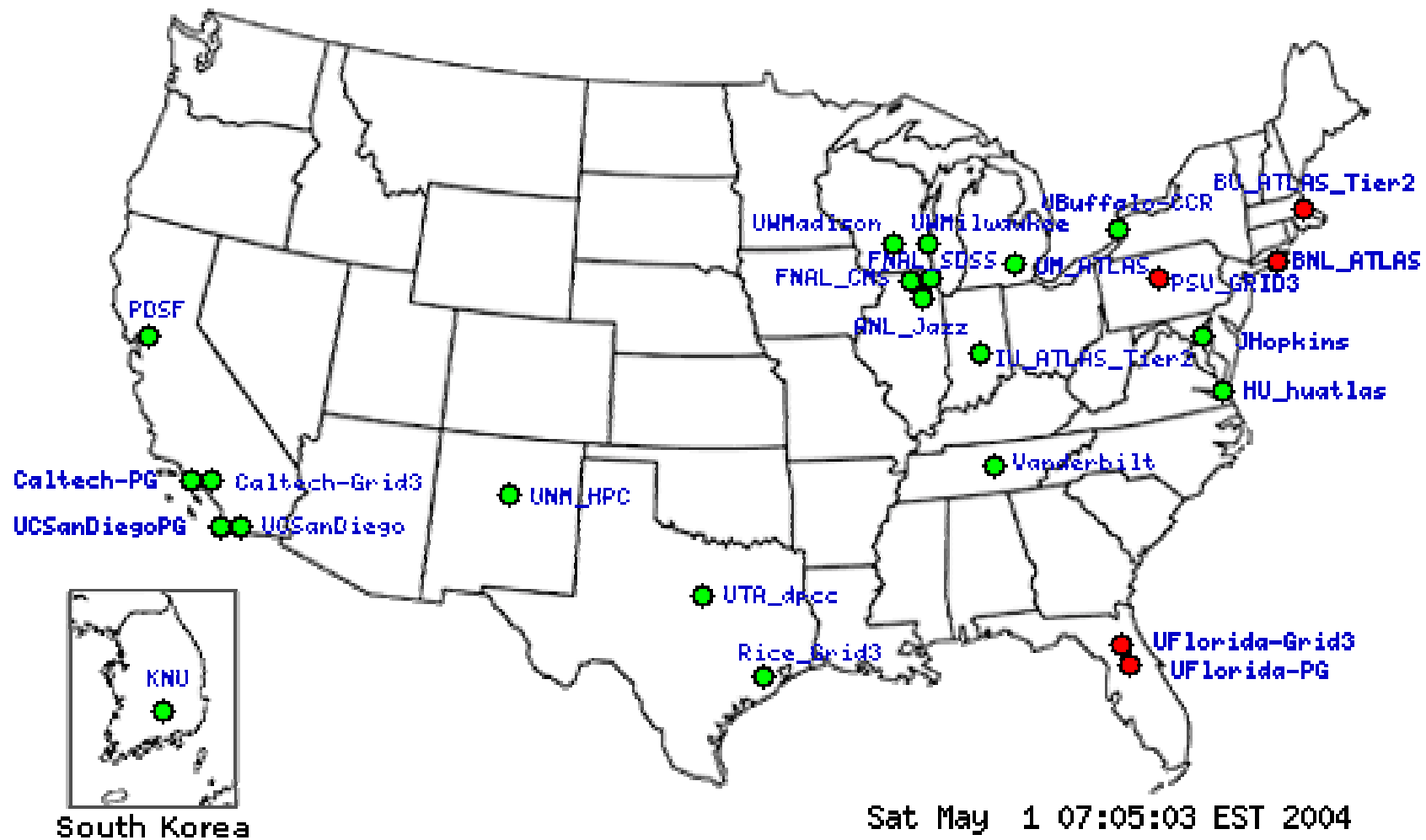


Organizations Managed as 4 Virtual-Organizations.

Applications including 4 from U.S. LHC.

Experiments including U.S. ATLAS and U.S. CMS Tier-1 & Tier-2

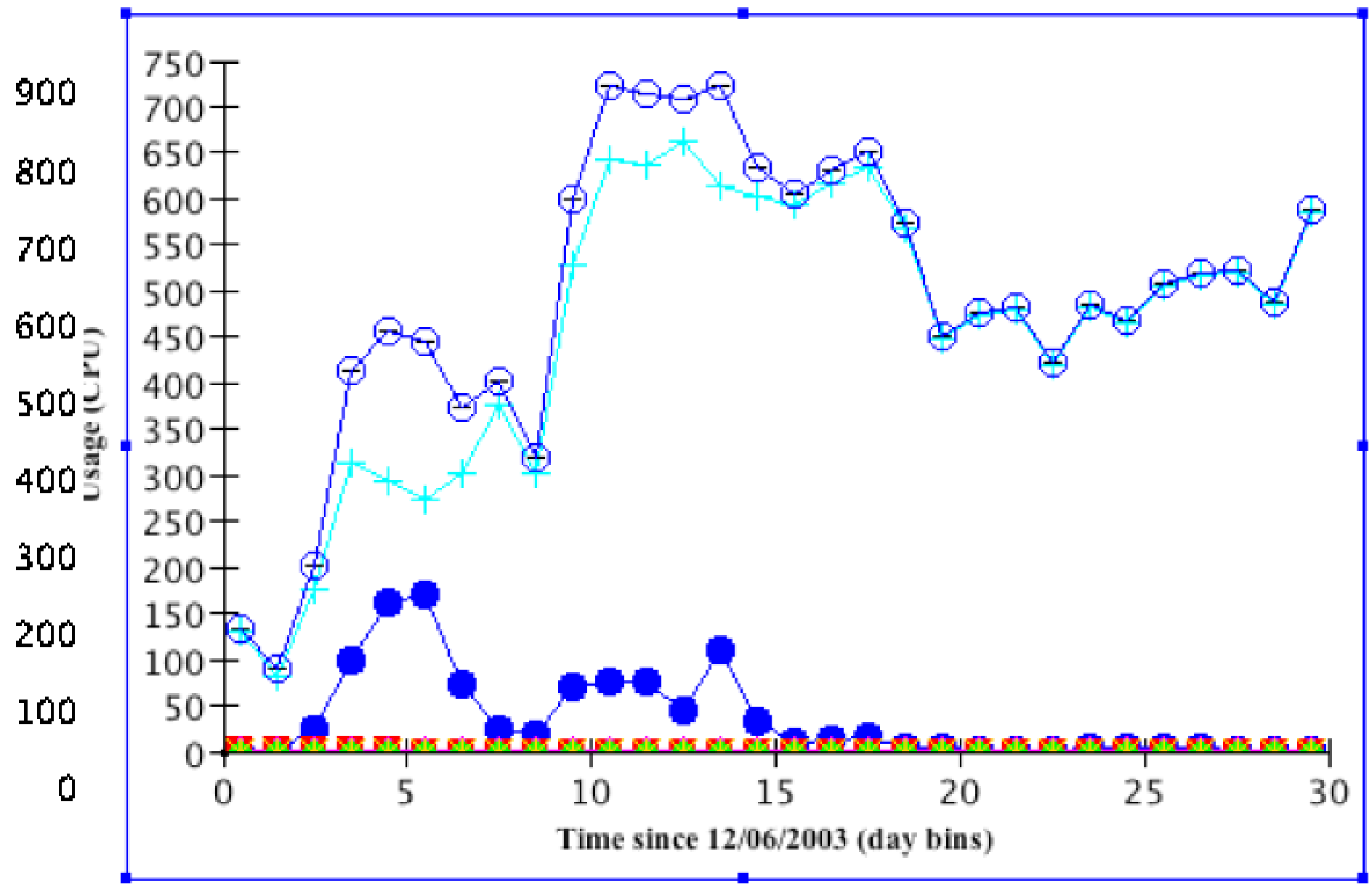
CPUs.



(site availability generated

by the site operators)

Usage: Running Jobs





Rob Gardner	Mike Wilde	Ian Fisk	Scott Koranda	Rich Baker
Jim Annis	Peter Couvares	Jorge Rodriguez	Leigh Grundhoeffter	Nickolai Kuropatine
Xin Zhao	John Hicks	Ed May	Alain Roy	Brian Moe
Fred Luerhing	Iowna Sakrejda	Yuri Smirnov	Marco Mambelli	Anzar Afaq
Suresh Singh	Carey Kireyev	Alain DeSmet	Jerry Gieraltowski	Doug Olson
Brian Tierney	Saul Youssef	Anne Heavey	Terrence Martin	Andrew Zahn
Scott Gose	Vijay Sekhri	Dantong Yu	Lawrence Sorrillo	Yong Xia
Rob Quick	Michael Ernst	Greg Graham	Bobby Brown	Bockjoo Kim
Jens Voekler	Ruth Pordes	Matt Allen	Yujun Wu	Lisa Giacchetti
Joe Kaiser	Erik Paulson	George Fekete	Dan Engh	Kihyeon Cho
James Letts	Tim Thomas	John Weigand	Iosif Legrand	Mark Green
Craig Prescott	Nosa Olomu	Ben Clifford	Dan Bradley	Timur Perelmutov
Patrick McGuigan	Shawn McKee	Guarang Mehta		

on in Grid3 was at the level of 58 people. 8 worked full time. 10 worked half  
strators worked quarter time.

ffort was at about the estimated 7 FTE-years (17 FTE equivalents for 5 mo

Metrics were defined in July and measured in November.

The Goal was to operate a Month. Organizations left their resources in Grid20

Performance Metrics were met, but not the Efficiency Metrics.

Metric	Goal	Achieved	Comments
Number of CPUs	500	2700	More than 60% of the CPUs are non-dedicated facilities. They are shared with local users
Number of Users	10	102	Most job submissions are by a small number of administrators.
Number of Sites	20	26	Complexity metric.
Number of Sites running Concurrent Applications	10	17	Demonstrates policies in use.
Data Transferred per day	2-3 TB	4 TB	GridFTP test application. Grid stable during test
Percentage of resources Used	90%	40-50%	% = Nodes used compared to those available to jobs. Note many nodes used by local jobs.
Efficiency of Job Completion	75%	75% and variable	
Number of Concurrent Jobs	1000	1100	< 50% of 2700. but many used by Local Jobs. are not measuring this metric correctly yet.
			Grid3 "Meltdown" due to revoked certificate (~3 times)

Architecture:

*ilities:* Execution and storage. Include non-dedicated, shared resources. No requirements on nodes.

*ices:* Processing, storage, account management, information, monitoring, configuration, operations.

*lications:* Installed dynamically without site administration. Application administrators responsible for the deployment, operation and monitoring of their applications.

Middleware:

Virtual Data Toolkit (VDT - Globus, Condor, NSF Middleware Initiative, EDG software extensions for VO (European Data Grid VO Management System), Information extension to Glue Schema) and monitoring (Globus MDS, U.S. CMS MonALISA, glia, MDViewer).

Each middleware (e.g. data management) is different for and the responsibility is different.

*g, Installation and Configuration:* Used iVDGL/ATLAS Pacman for packaging and installation of middleware and applications. Goal to make installation simple and minimize

## ing Service

Condor-G or Chimera Virtual Data System submissions through Globus-Gram keeper to one of 4 standard batch systems (PBS, Condor, LSF, FBSng).

## nsfer Services

FTP interfaces on all sites through gateway systems

Files are transferred into processing sites using globus-url-copy.

Application specific transfer of Results back

- Basic: GridFTP based Data Movement to Permanent Storage System
- Advanced: Managed Data Storage and Data Access Services based on SRM/dCache

## ation and Accounts

Certificates; Globus Gridmapfiles; Unix Account for Each VO.

## g Services

System and application level monitoring

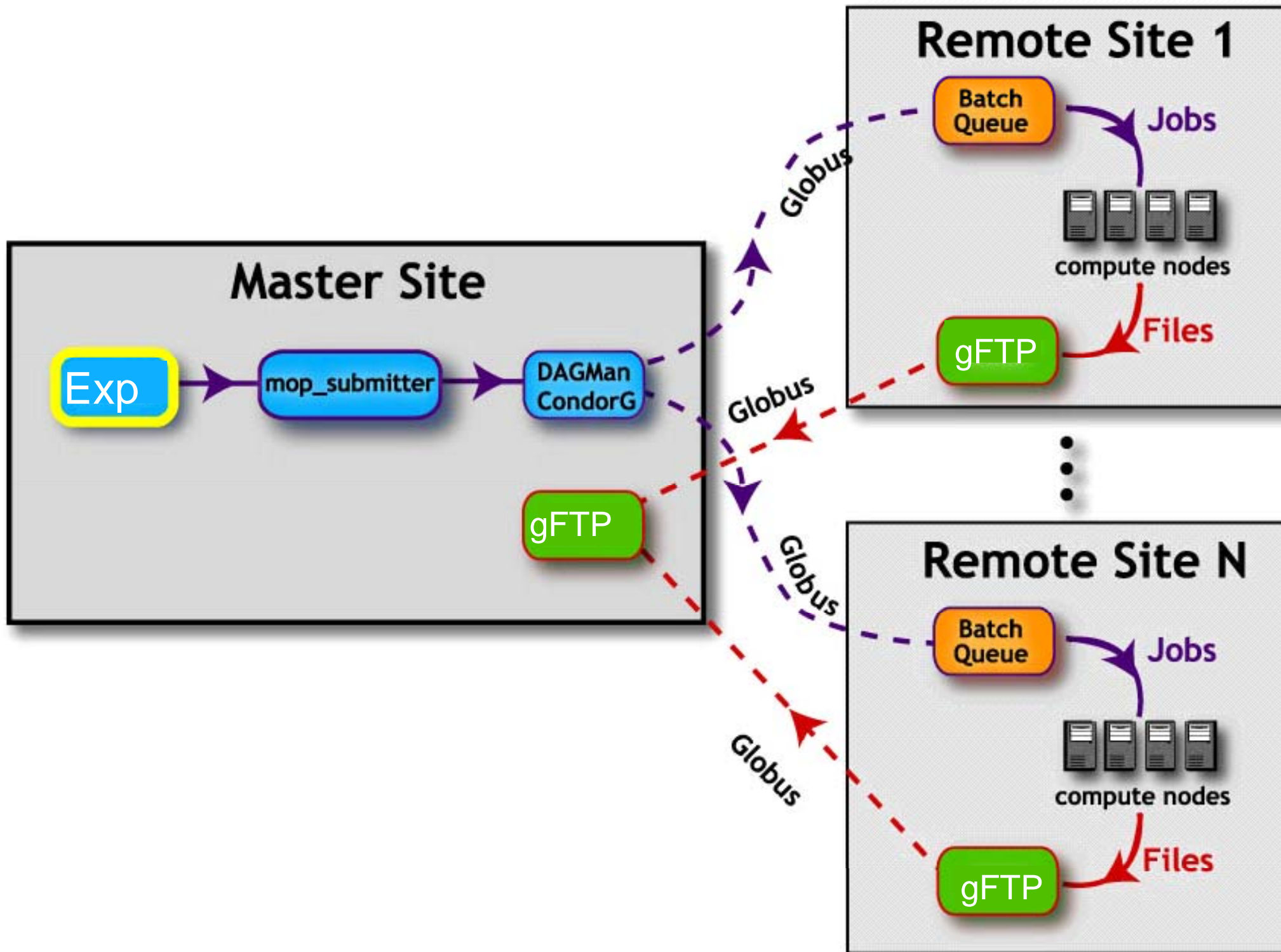
Integrators between different implementations

## rganization Management Services

WebS server + administration interface per VO.

Scripts to automatically generate gridmapfiles.

## on Services



Grid "Regional Center" used dedicated US CMS resources for

releasing Grid3 resources for CMS Simulation Production

**100% Gain In**

**CPU for CMS**

**through Grid**



CMS Production Facilities stay part of Grid3

CMS simulation production fully compatible with Grid3 environment  
All essential US CMS required functionalities have “production quality” support

opportunistic use of non-CMS resources continues to be successful!

After SC2003 milestone, Grid3 kept running successfully  
Some local configuration issues, did not impact overall stability  
Looks like a successful strategy (effectively doubling resources)

Integration of U.S. Grid with the rest of LHC Grid?

Operations model that would support development and new releases  
Technology cycles (data analysis!) on Grid3 and LCG/EGEE?



has Demonstrated:

heterogeneous facilities can be used in a common Grid Environment.

operate and use a Multi-Organization Grid with distributed ownership and a coherent system.

Grid of over ~20 Facilities can be robust and performant for simple production applications.

feasibility of the strategy of federating and sharing resources - Open Science Map.

Plans:

rate the current infrastructure for Data Challenges

live the common Grid Environment for increased capability and performance

t longer term Engineering of Services and Capabilities aligning with and contributing to the Open Science Grid.

continue interoperability and joint projects with the LCG



stitutions participating in Grid2003 will continue to contribute their resources and expertise to the shared common infrastructure in accordance with the MOU with LBNL. The details of a model for ongoing Operations will be worked out over the next few months.

Grid2003+ will Operate and Evolve the current Grid Environment:

- Upgrade infrastructure to new versions of middleware and applications;
- Follow up on recommendations from Lessons Learned document (robustness and hardening, extended monitoring, operations infrastructure)
- Continue adding and integrating a (SRM/dCache based) Storage Element.

Grid2003+ will collaborate on further projects to Engineer and Deploy the Next Phase of Grid Computing resources and technologies

Grid2003+ is part of the continuation of the existing Grid projects in the US and LCG and will be done in conjunction with the Open Science Grid engineering & blueprint

ing strategy of Interoperability and Joint Projects with the LHC  
outing Grid Project on many fronts.

orative Efforts in particular on

ommon Virtual Data Toolkit (VDT) delivery and support team.

ata Movement and Storage Management: U.S. CMS demonstration  
etween Grid2003 and Cern. Using GridFTP, dCache, Storage Resource  
anagement (SRM).

ob Execution: U.S. ATLAS Grid3 application submission to LCG sites using  
himera Virtual Data System (VDS).

erge of Information Attributes (GLUE Schema extensions) from Grid2003  
nd LCG.

### Collaborative Efforts:

rtual Organization Management Project (VOX) collaboration with European  
ata Grid and LCG Security Working Group.

ontributions to and from the wider CMS and ATLAS software and  
omputing deliverables.

resentations at and discussions with LCG committees (Grid Deployment  
oard, Project Oversight Board, Software Computing Committee, Project  
xecution Board)

articipation in High Energy Physics Joint Technical Board and Global Grid  
orum Particle and Nuclear Physics Applications Research Group



Putting “real”, multi-organizational Grids to work

not talk about EGEE

the HEP Collaboration (or parts of it) fundamentally is a “physicists organization”

R&D project with an operational component

not a sustained IT organization or a “distributed computing center”

need to build the partnerships and need to address the organizational issues

how to build the supporting structures to run a truly distributed, engineered

managed, robust and resilient, accountable and secure service for data access

and analysis on Peta-scales!

working with centers, universities and projects to formulate a

roadmap towards the “Open Science Grid”

LHC as an exemplar global science project to drive the creation of a US National

resource for science – the Open Science Grid

US grid middleware and services the basis for international initiatives to build

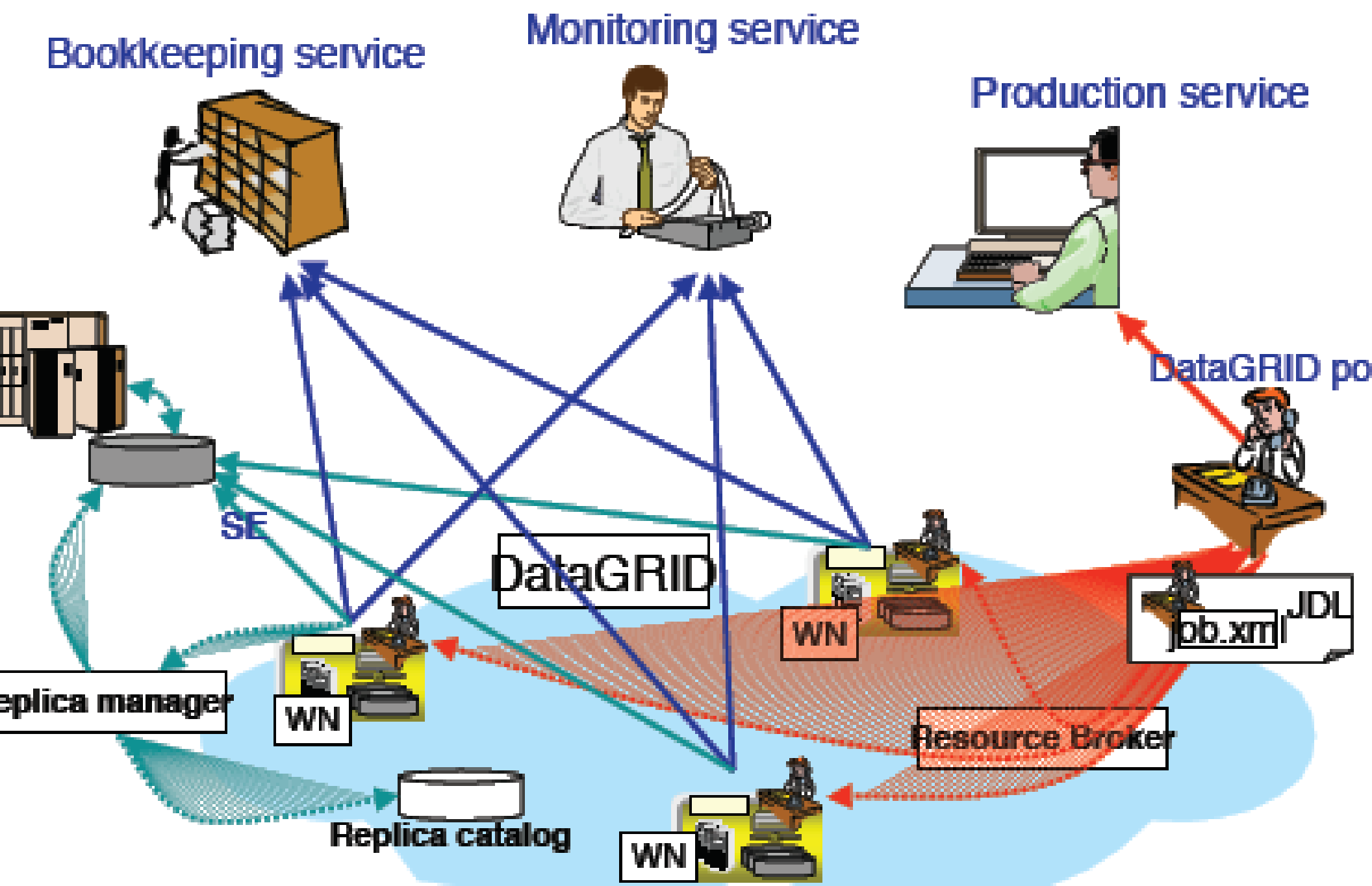
production grids for science

it is time to federate the (U.S.) resources and to continue to “lead”

in constructing a global grid

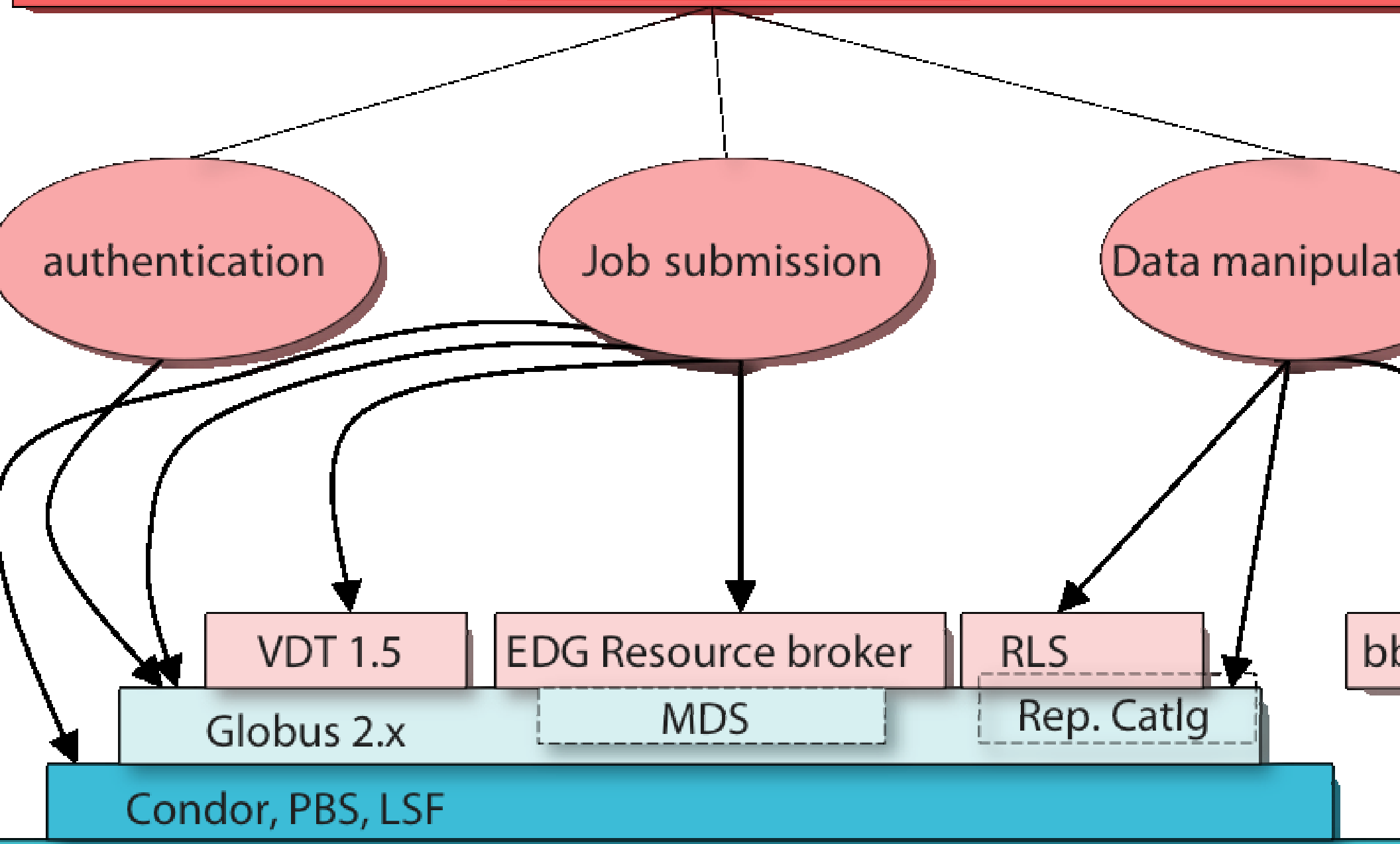
to build the OSG so it is open to other sciences and complements and interoperates

services for processing, to be further developed



## LHC software stack for Event Simulation Production

### LHC Event Simulation Production Use Cases



# Layered Grid Architecture

(I. Foster et al.)

Architecture: (H. Newman)

Above the Collective Layer

Application Codes

Reconstruction, Calibration, Analysis

Users' Software Framework Layer

Modular and Grid-aware:

Architecture able to interact effectively

with the lower layers (above)

Applications Layer

Users and algorithms that govern system operations)

Policy and priority metrics

Workflow evaluation metrics

Task-Site Coupling proximity metrics

End-to-End System Services Layer

Workflow monitoring and evaluation mechanisms

Error recovery and long-term redirection mechanisms

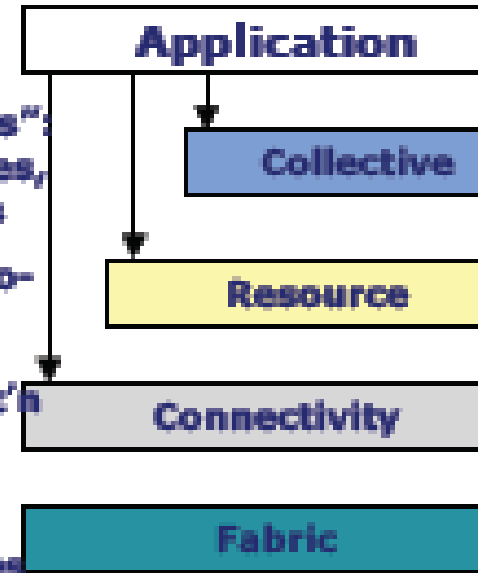
System self-monitoring, steering, evaluation and optimization mechanisms

"Coordinating multiple resources":  
ubiquitous infrastructure services,  
app-specific distributed services

"Sharing single resources": nego-  
tiating access, controlling use

"Talking to things": communicat'n  
(Internet protocols) & security

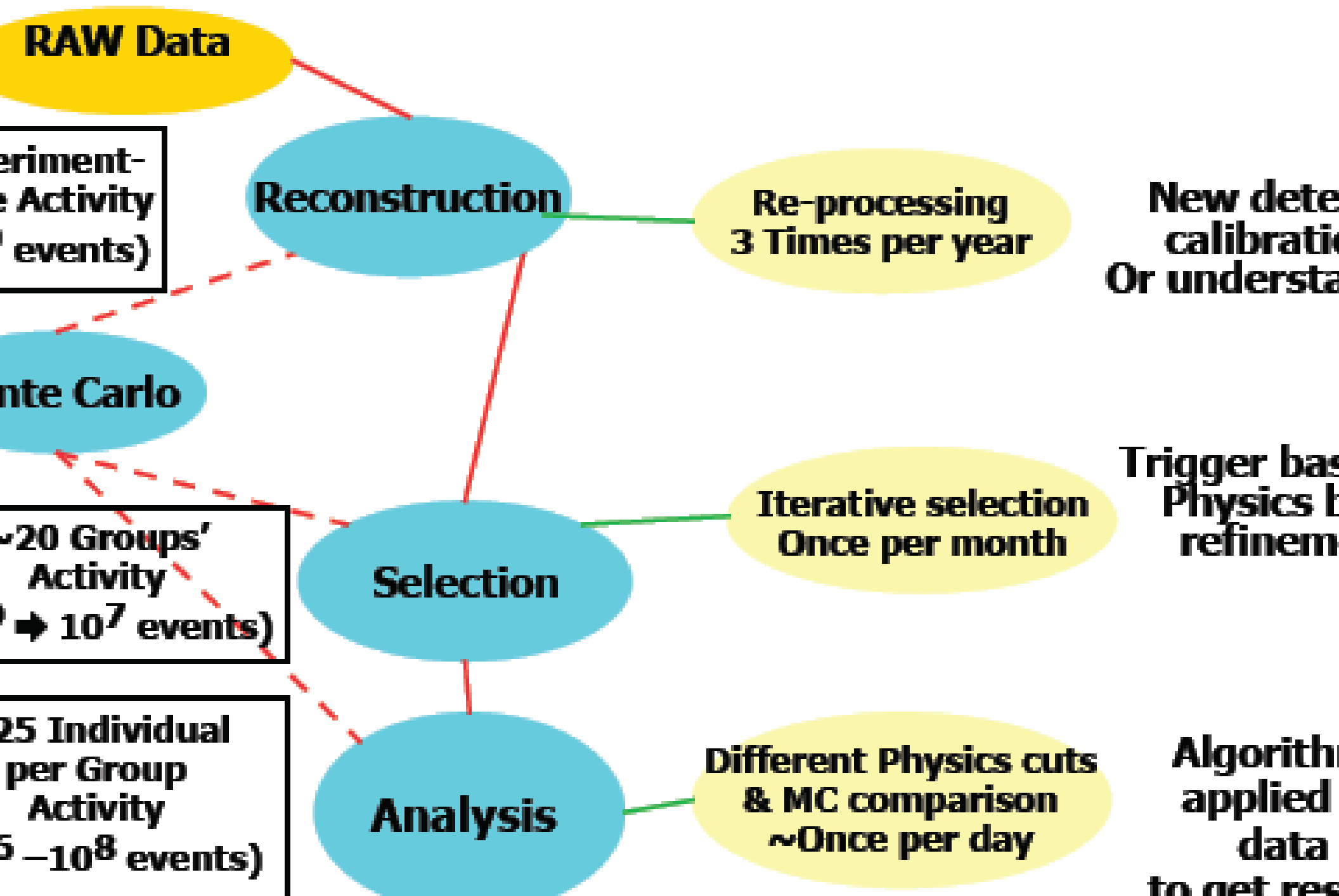
"Controlling things locally":  
Access to, & control of, resources



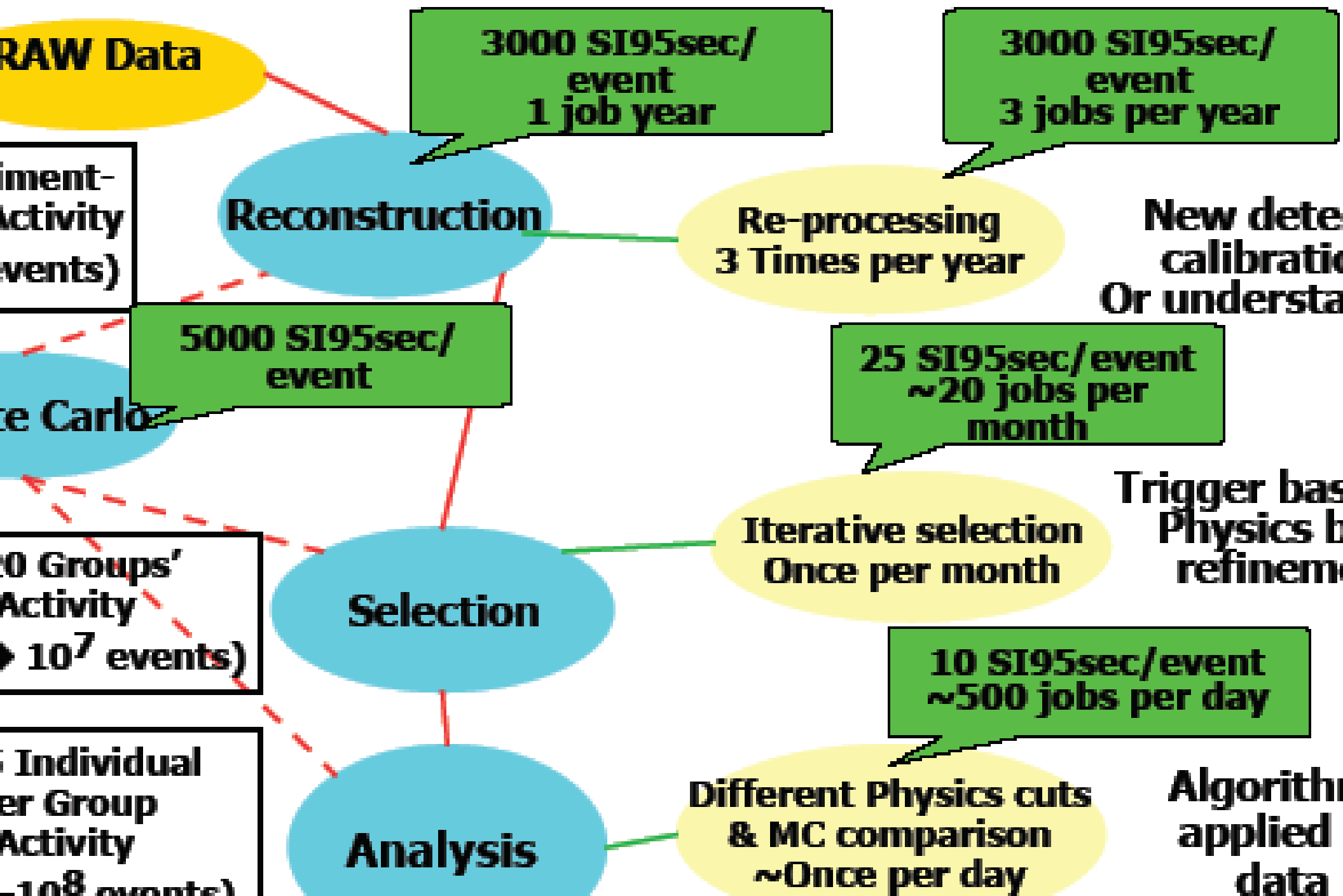
forward to physics data analysis means a significant paradigm shift  
from well-defined production jobs => interactive user analysis  
from DAGs of process => “Sessions” and state-full environments  
from producing “sets of files” => accessing massive amounts of data  
from files => data sets and collection of objects  
from using essentially “raw data” => complex layers of event representation  
from “assignments” from a central repository => Grid-wide queries  
from “user registration” => enabling sharing and building communities  
Are (Grid) technologies ready for this?  
There needs to be a tight inter-play between prototyping the analysis services  
developing the “lower level” services and interfaces => ARDA Prototype  
These are going to be the “new paradigms” that will be exposed to the user  
Can a “data analysis session” transparently extended to a distributed system?  
Yes, but requires a more prescriptive and declarative approach to analysis  
What services for “collaborative” work?  
What new paradigms beyond “analysis”



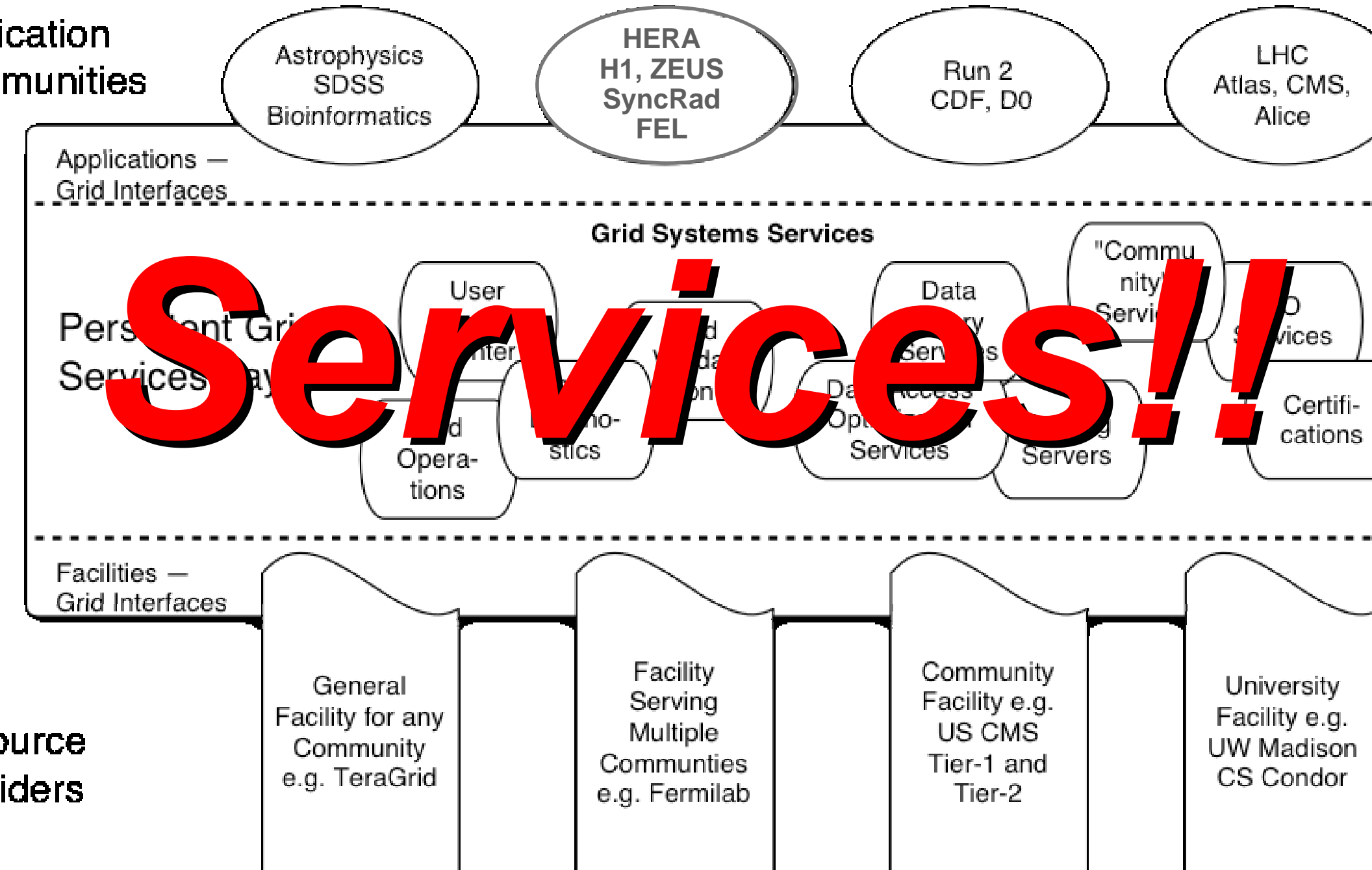
# Hierarchy of Processes (Experiment, Analysis Groups, Individuals)



# Hierarchy of Processes (Experiment, Analysis Groups, Individuals)



# Grid Layer “Abstraction” of Facilities — Rich with Services!



m extracts a subset of the datasets from the virtual file catalogue  
data conditions provided by the user.

m splits the tasks according to the location of data sets.

balancing between local data access and data replication.

n sub-jobs and submit to Workload Management with precise job  
ptions

User can control the results while and after data are processed

t and Merge available results from all terminated sub-jobs on rec

sis objects associated with the analysis task remains persistent i

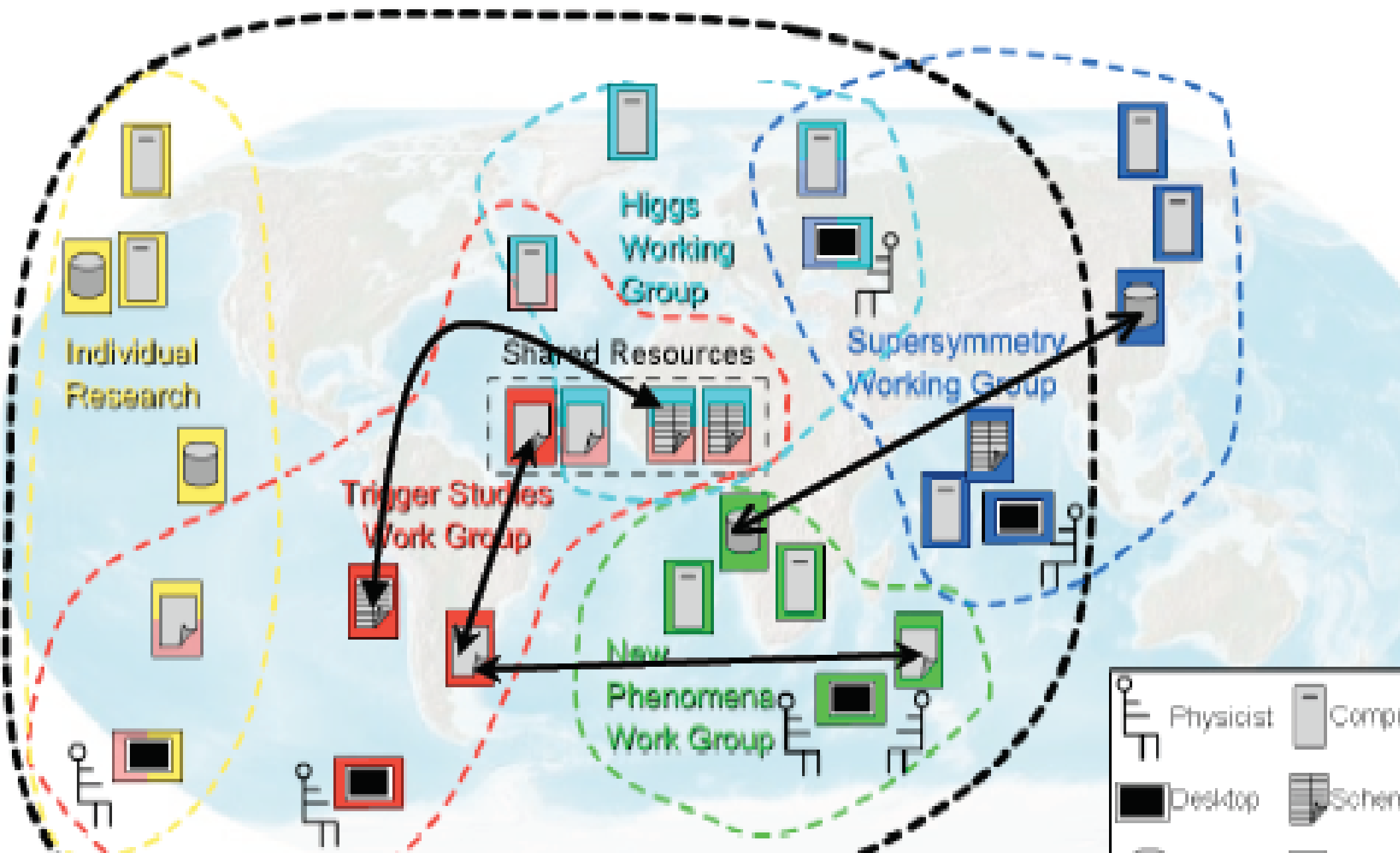
nvironment so the user can go offline and reload an analysis tas

late, check the status, merge current results or resubmit the sam

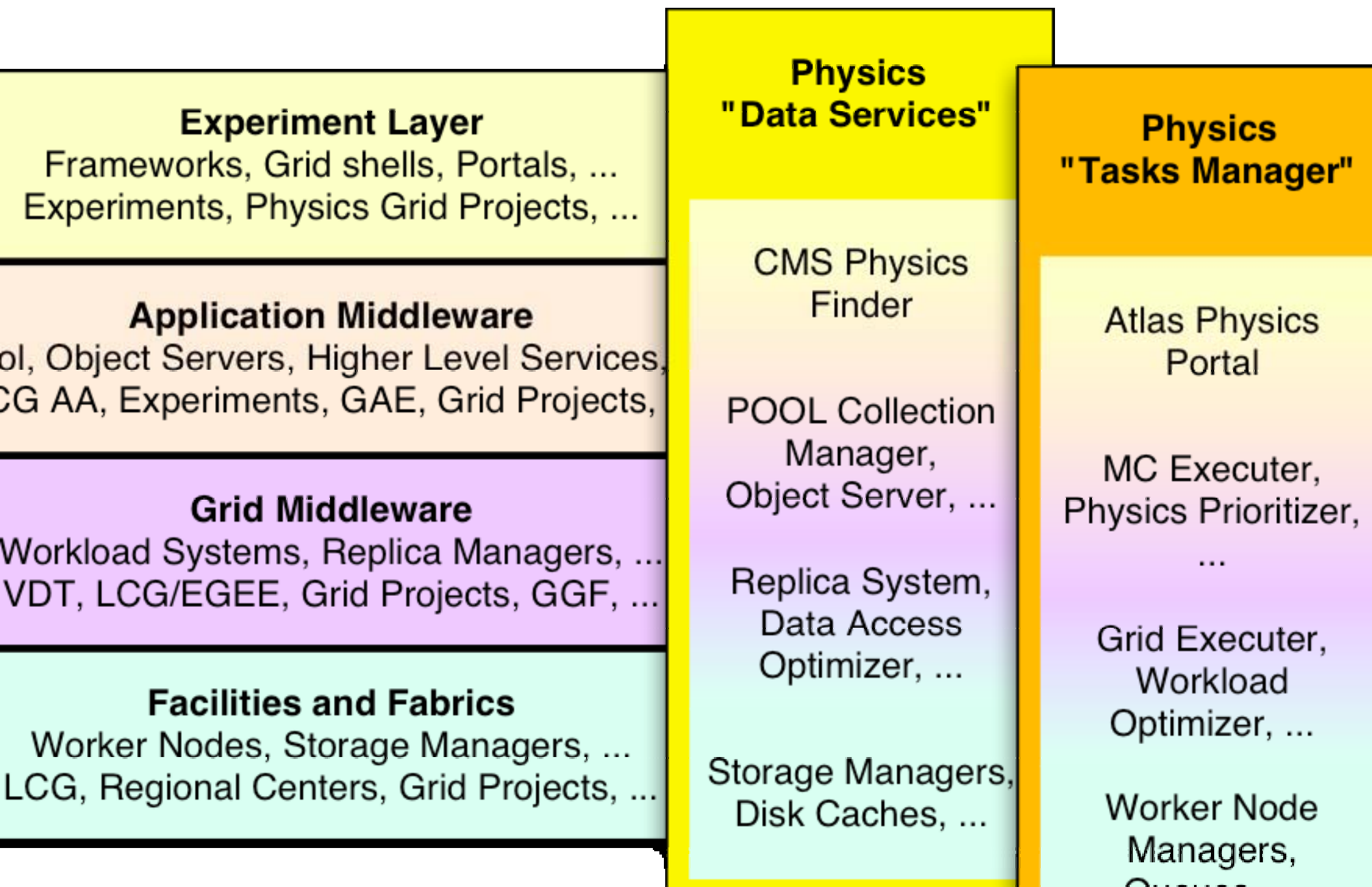
modified analysis code.

# Communities of Scientists Using the Grid for Distributed Analysis

Infrastructure for sharing, consistency of physics and calibration data, software



# Experiment's "Services" go end-to-end!



# What is analysis anyway?

Monday, 10 May, 2004

Massimo Lamanna (CERN)

“The ARDA Project: Grid Analysis Prototypes  
of the LHC Experiments”

's HEP Collaborations are getting ready for the big challenges  
expected to lead to discoveries of new elementary particles and new  
behaviors of the fundamental forces — and a decades-long scientific  
program  
adopted a globally distributed computing model, to enable science  
at global scale, and to enable scientists worldwide to be full part of the  
expected breakthrough discoveries  
many technology and organizational challenges ahead  
collaborating with computer scientists and other scientific  
communities  
construct a global cyberinfrastructure of international computing grids  
that will be used by thousands of scientists  
goal to enable scientific collaborators to work together as co-located  
peers, and to create new capabilities to empower the individual scientists



# of Research for Grid Federations

Interface Languages and Standard Protocols

Experiments need to be able to describe Interfaces needed by applications

Current model where Computing Resource Providers encourage to install suites of software limits ability to flexibly deploy Grid Services for Applications

Monitoring and Information Providing

to track VO usage of resources

to indicate to incoming users what the likely priority is they will receive

much richer information provider is needed to convey information to optimize

users to enable intelligent scheduling decisions

Authentication – A lot of work has been done, but ...

Consensus needs to be built on how tools are used and administered

Authorization and Privilege – Almost nothing exists

Closely related to priority and quota setting

erating Resources to form a Federated Grid is necessary  
g together Computing Resources from Different Commu  
only then that the Power of Grid Enabled Distributed Co  
begin to be realized.

important for the Computing Providers and the Experime  
ed to work within the Framework of Federated Resource  
beginning of Development.

veloping a Homogeneous Distributed Computing System  
hes the Experiments very little about Operating in a Grid  
ronment, and more importantly gains them very little in te  
puting Resources.