



DØ Computing Experience and Plans for SAM-Grid

EU DataGrid

Internal Project Conference

May 12-15, 2003

Barcelona

Lee Lueking

Fermilab

Computing Division

Roadmap of Talk

- DØ overview
- Computing Architecture
- SAM at DØ
- SAM-Grid
- Regional Computing Strategy
- Summary



May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.



The DØ Experiment

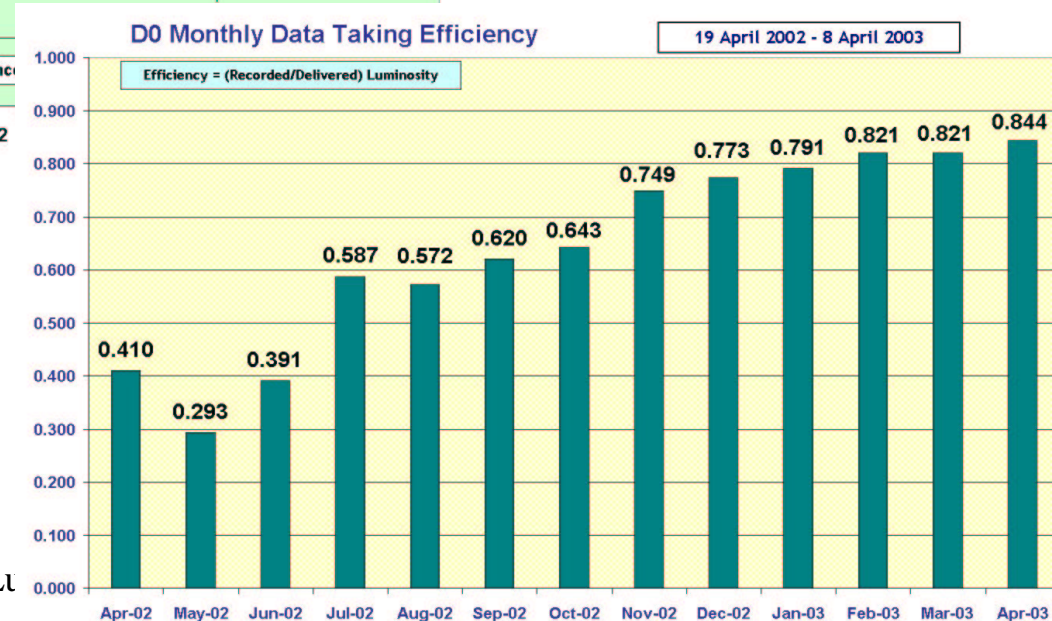
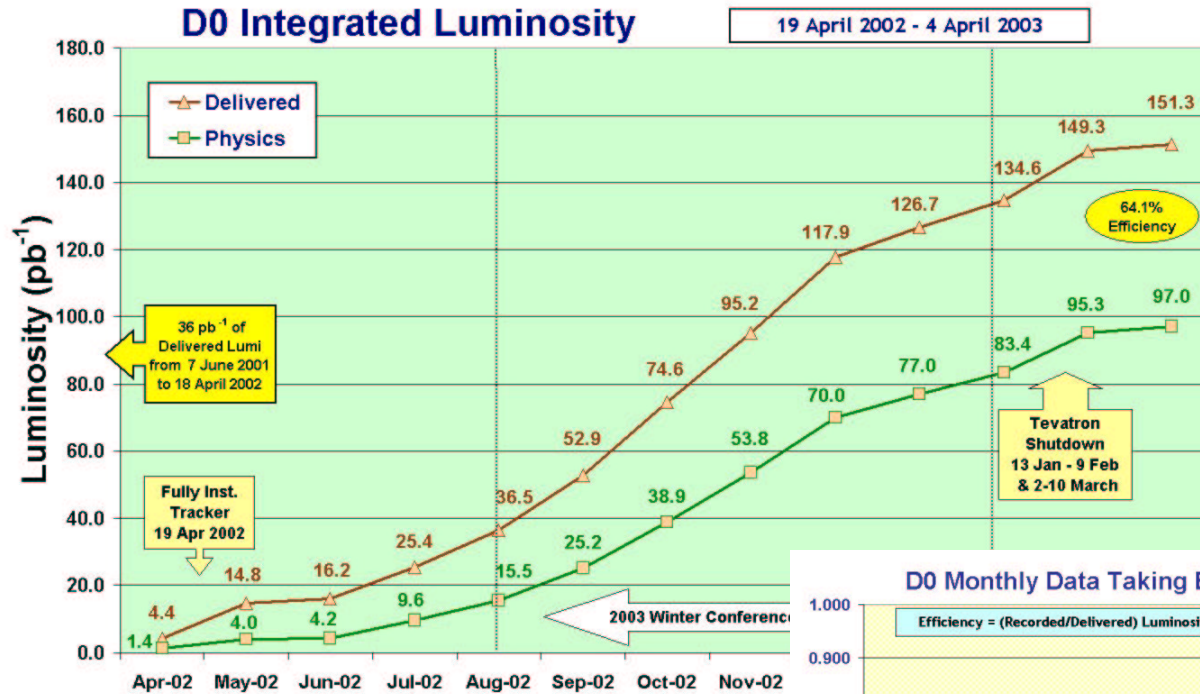


- DØ Collaboration
 - 18 Countries; 80 institutions
 - >600 Physicists
- Detector Data (Run 2a end mid '04)
 - 1,000,000 Channels
 - Event size 250KB
 - Event rate 25 Hz avg.
 - Est. 2 year data totals (incl. Processing and analysis): 1×10^9 events, ~1.2 PB
- Monte Carlo Data (Run 2a)
 - 6 remote processing centers
 - Estimate ~0.3 PB.
- Run 2b, starting 2005: >1PB/year





DØ Experiment Progress



May 12-15, 2003

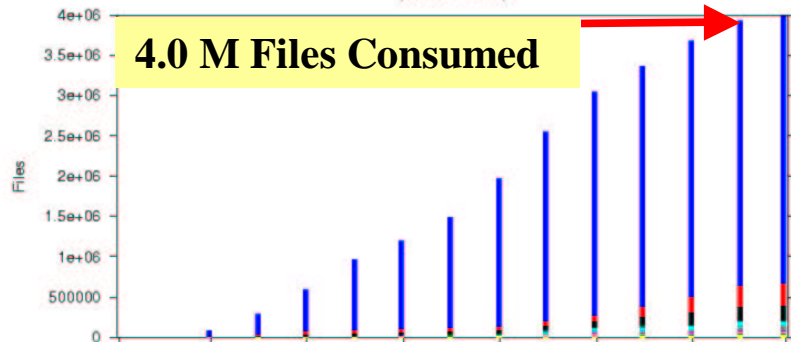
Lee Lu



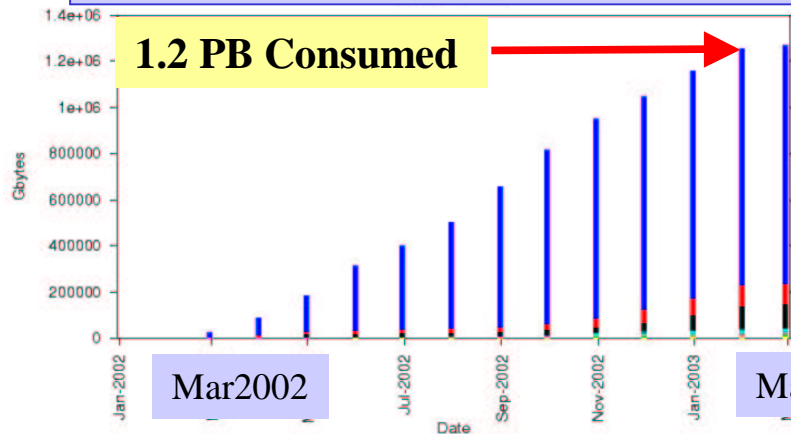
Overview of DØ Data Handling



Integrated Files Consumed vs Month (DØ)

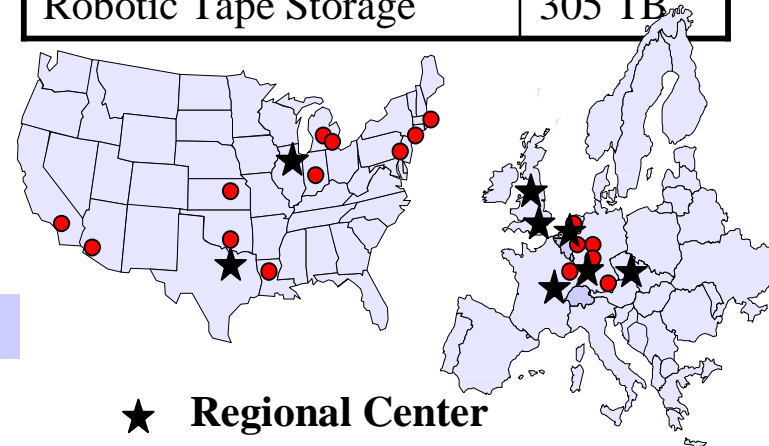


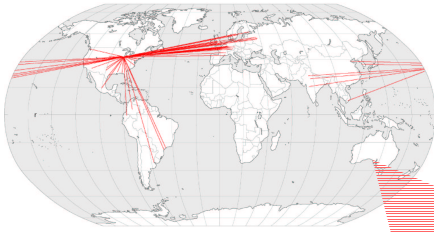
Integrated GB Consumed vs Month (DØ)



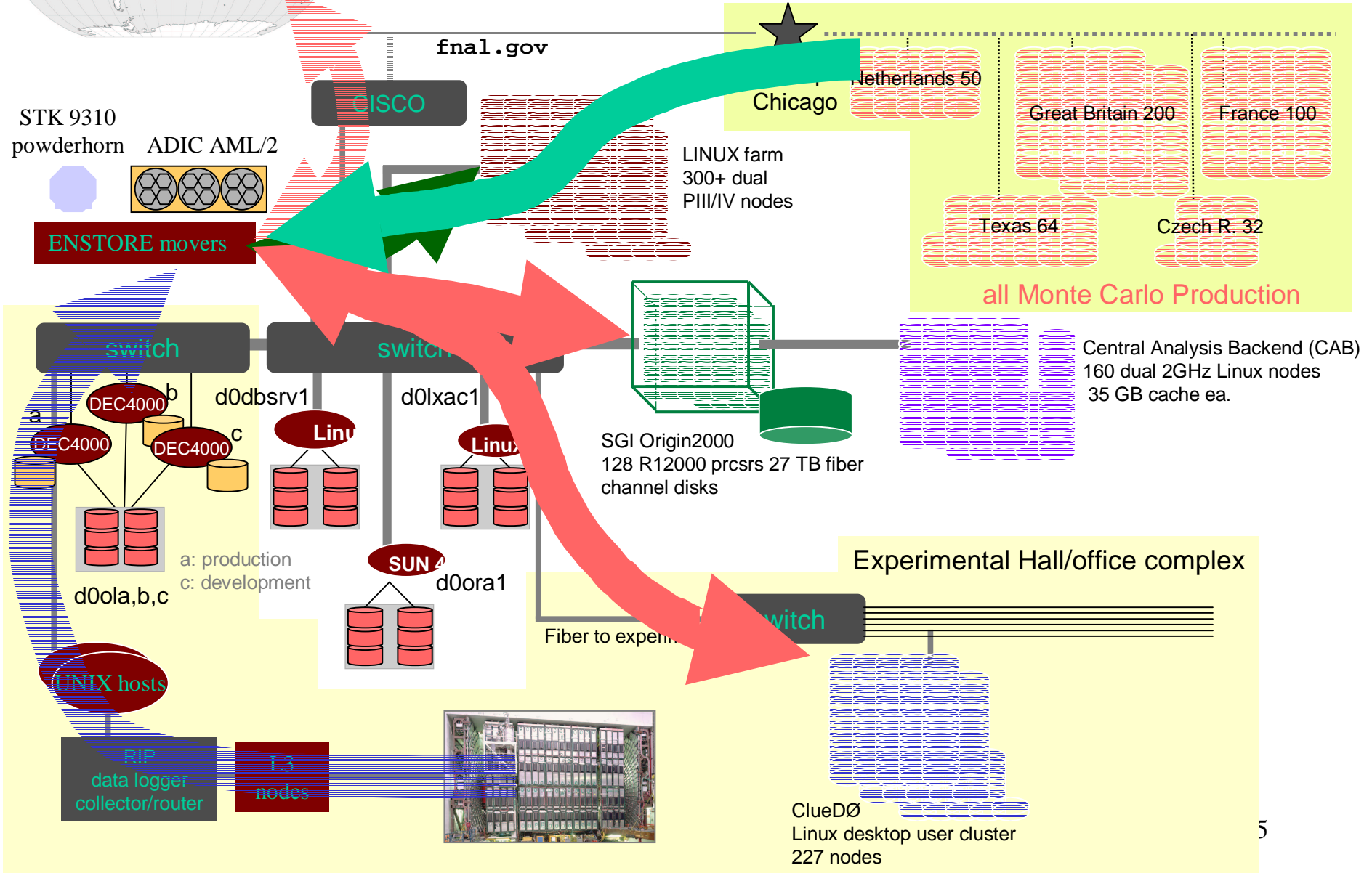
Summary of DØ Data Handling

Registered Users	600
Number of SAM Stations	56
Registered Nodes	900
Total Disk Cache	40 TB
Number Files - physical	1.2M
Number Files - virtual	0.5M
Robotic Tape Storage	305 TB





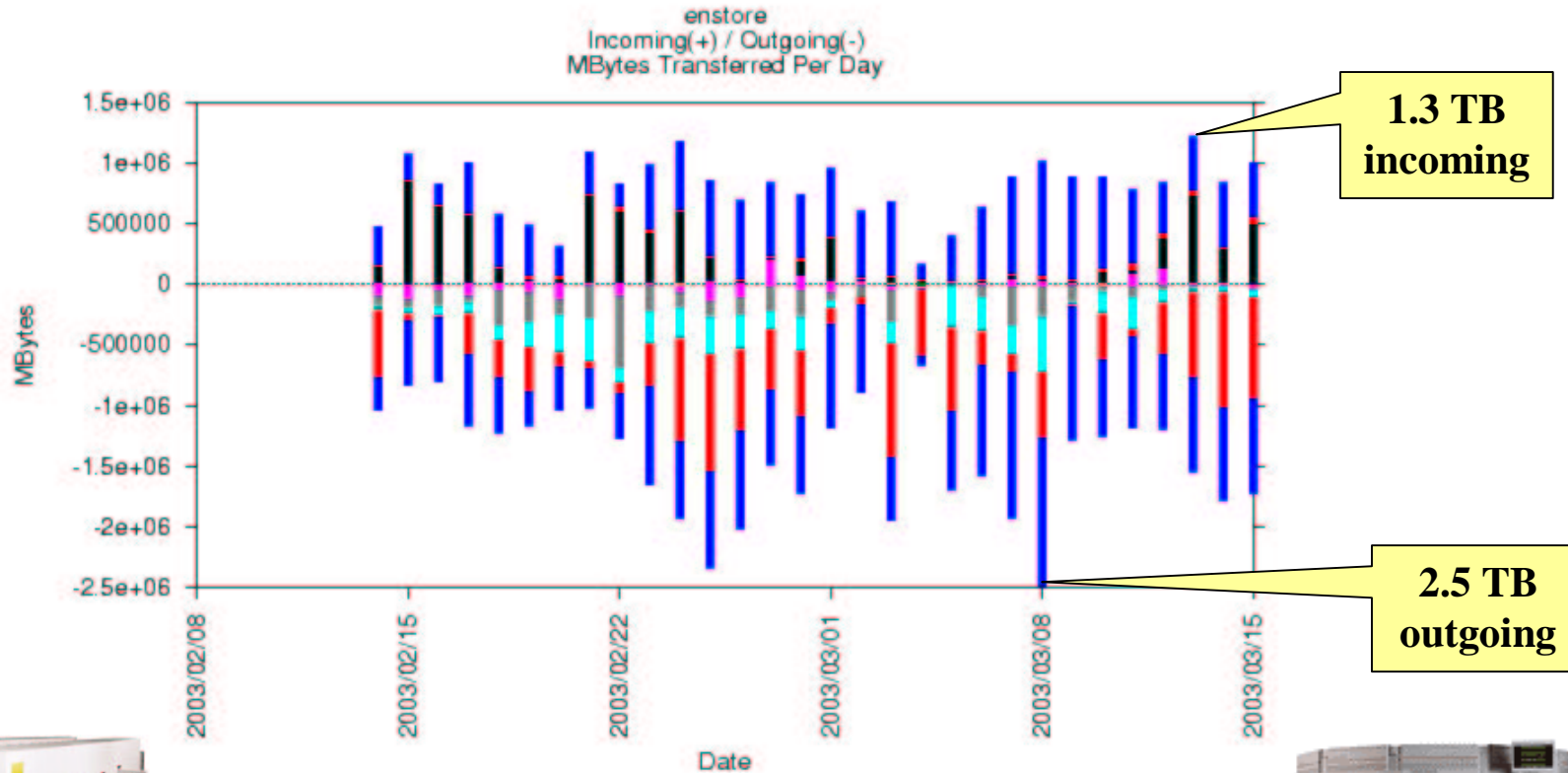
DØ computing/data handling/database architecture





Data In and out of Enstore

(robotic tape storage) Daily Feb 14 to Mar 15



Stations:

- fnal-farm
- ral-analysis
- datalogger-d0olc
- central-router
- cab
- cloud0
- umdzero
- other



May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.



SAM at DØ



d0db.fnal.gov/sam



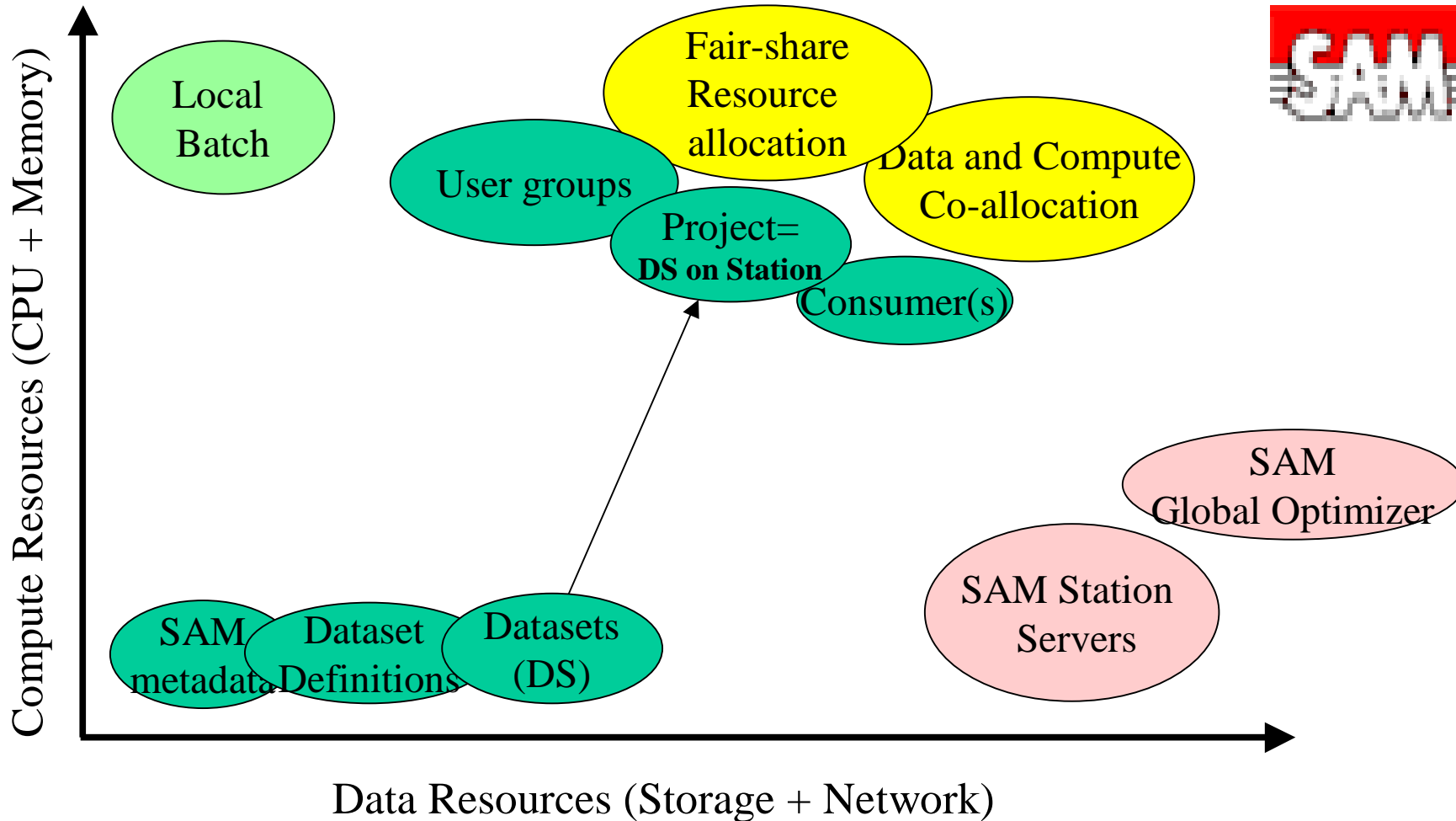
May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.

7



Managing Resources in SAM



● Batch scheduler
 ● SAM Meta-data
 ● SAM servers
 ● Batch + SAM

May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.



SAM Features



- **Flexible and scalable model**
- **Field hardened code**
- **Reliable and Fault Tolerant**
- **Adapters for many local batch systems: LSF, PBS, Condor, FBS**
- **Adapters for mass storage systems: Enstore (FNAL), HPSS (Lyon), and TSM (GridKa)**
- **Adapters for Transfer Protocols: cp, rcp, scp, encp, bbftp, GridFTP**
- **Useful in many cluster computing environments: SMP w/ compute servers, Desktop, private network (PN), NFS shared disk,...**
- **User interfaces for storing, accessing, and logically organizing data**





DØ SAM Station Summary



Name	Location	Nodes/cpu	Cache	Use/comments
Central-analysis	FNAL	128 SMP*, SGI Origin 2000	14 TB	Analysis & DØ code development
CAB (CA Backend)	FNAL	16 dual 1 GHz + 160 dual 1.8 GHz	6.2 TB	Analysis and general purpose
FNAL-Farm	FNAL	100 dual 0.5-1.0 GHz +240 dual 1.8 GHz	3.2 TB	Reconstruction
CLueDØ	FNAL	50 mixed PIII, AMD. (may grow >200)	2 TB	User desktop, General analysis
DØkarlsruhe (GridKa)	Karlsruhe, Germany	1 dual 1.3 GHz gateway, >160 dual PIII & Xeon	3 TB NFS shared	General/Workers on PN. Shared facility
DØumich (NPACI)	U Mich. Ann Arbor	1 dual 1.8 GHz gateway, 100 x dual AMD XP 1800	1 TB NFS shared	Re-reconstruction. workers on PN. Shared facility
Many Others > 4 dozen	Worldwide	Mostly dual PIII, Xeon, and AMD XP		MC production, gen. analysis, testing



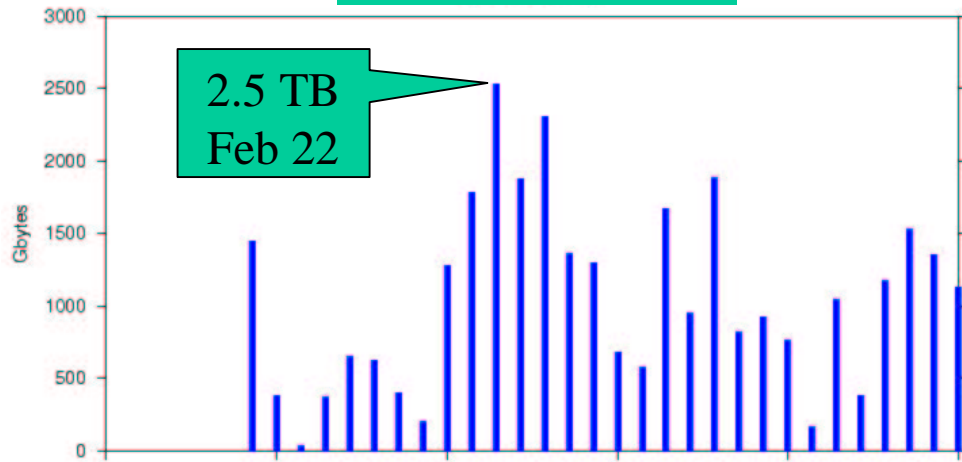


Station Stats: GB Consumed

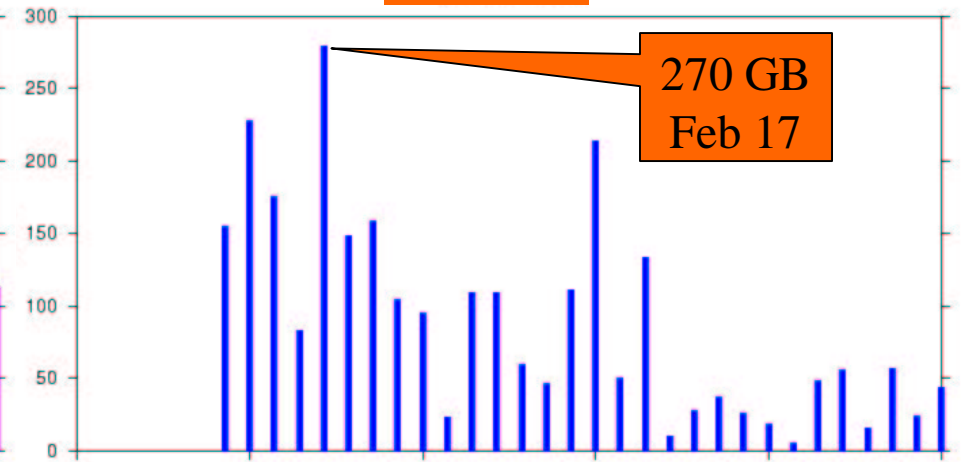
Daily Feb 14 – Mar 15



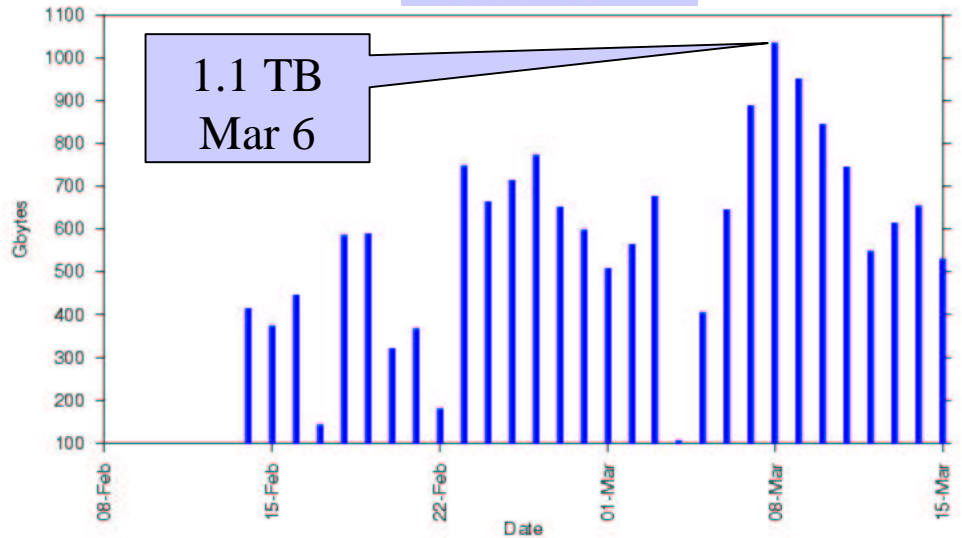
Central-Analysis



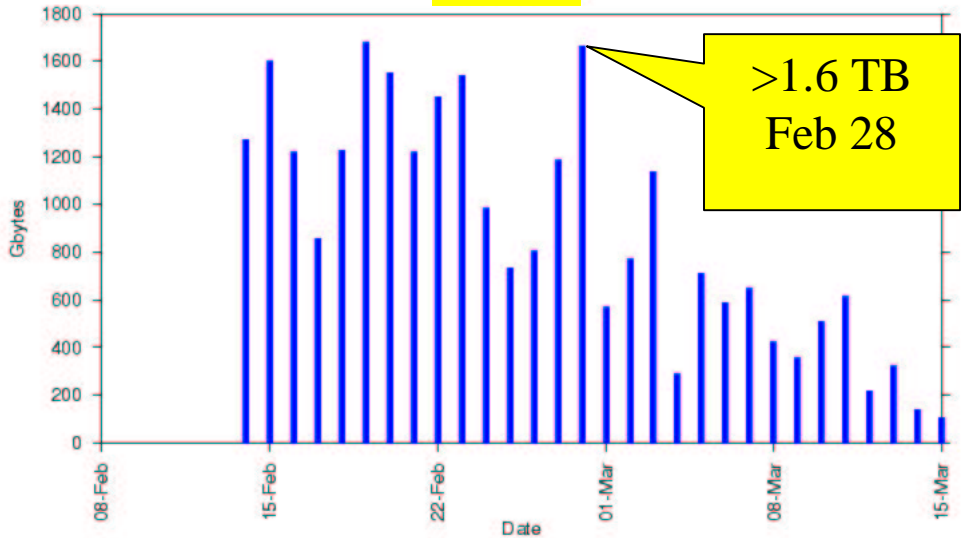
ClueD0



FNAL-farm



CAB



Station fnal-farm

Station cab

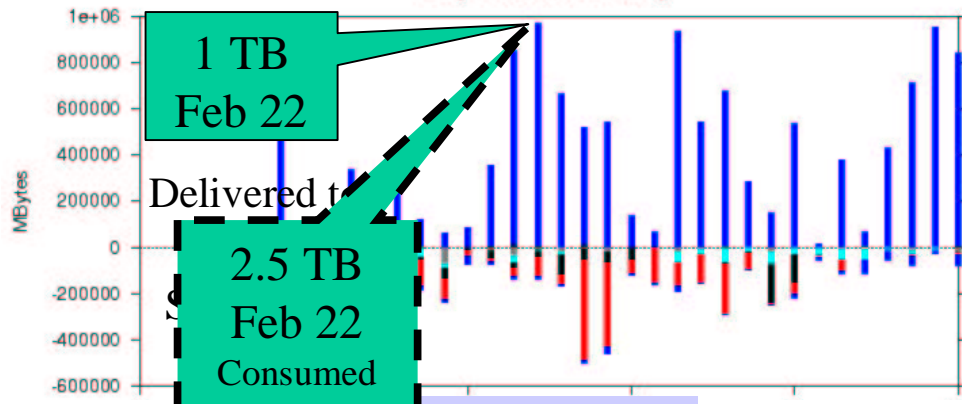


Station Stats: MB Delivered/Sent

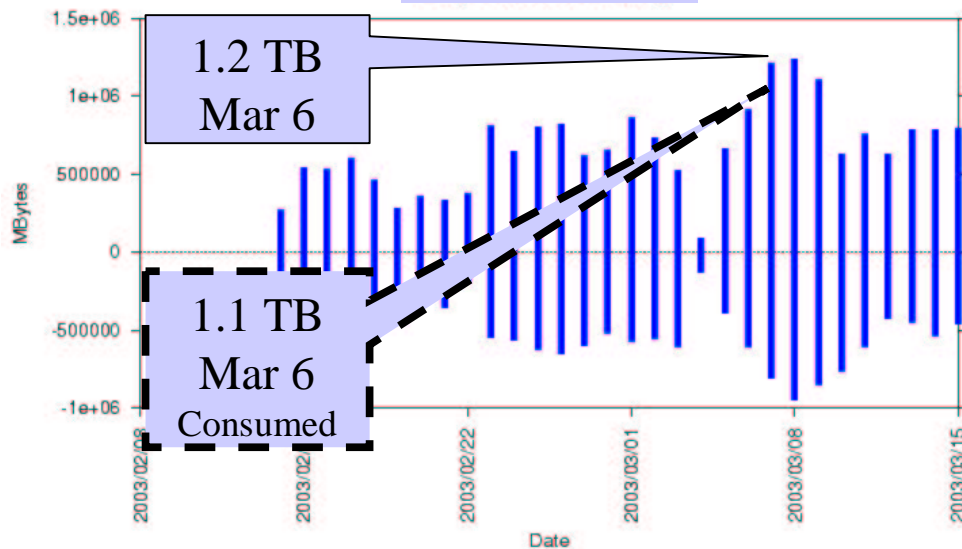
Daily Feb 14 – March 15



Central-Analysis



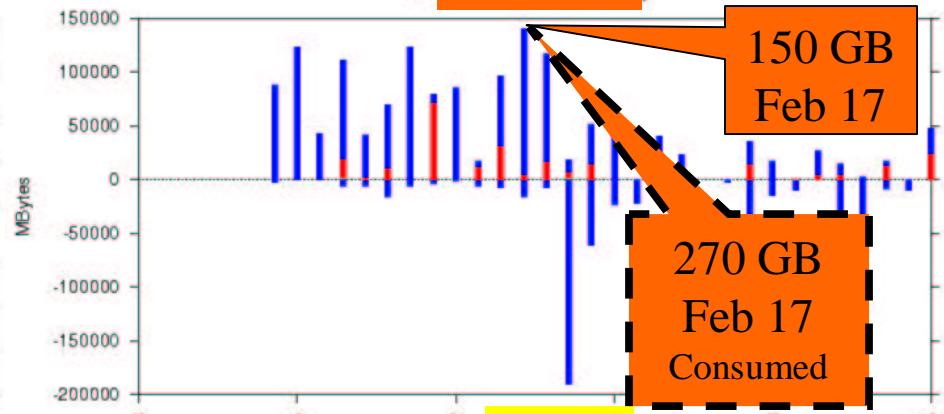
FNAL-farm



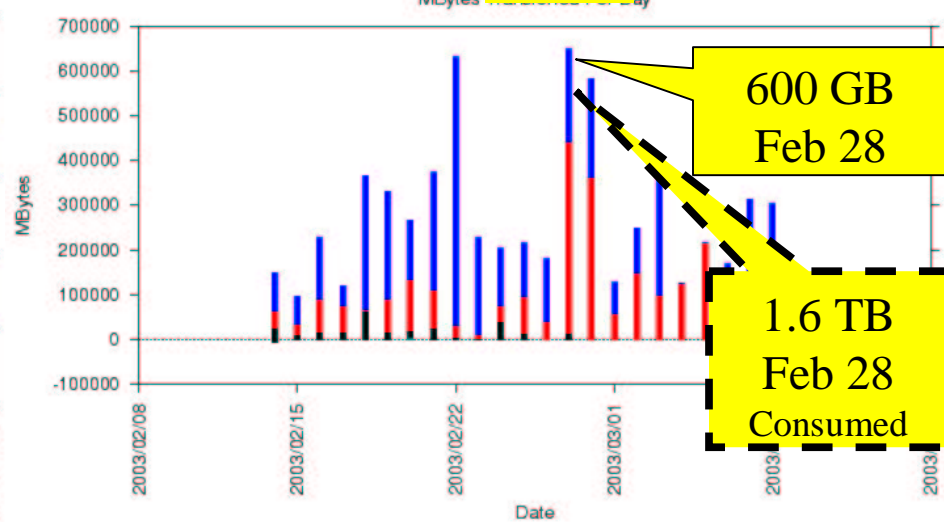
Stations:
enstore

clued0

ClueD0



CAB



Stations:
 enstore central-router uta-hep d0karlsruhe
 central-analysis princeton-d0 imperial-test umslara

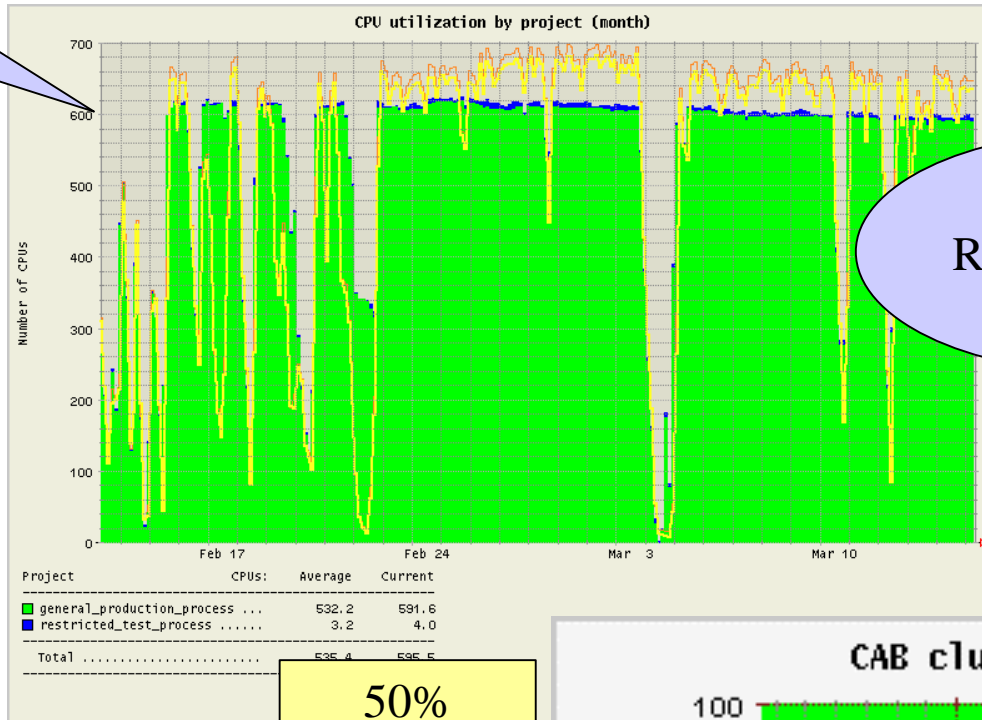


FNAL-farm Station and CAB CPU Utilization



Feb 14 – March 15

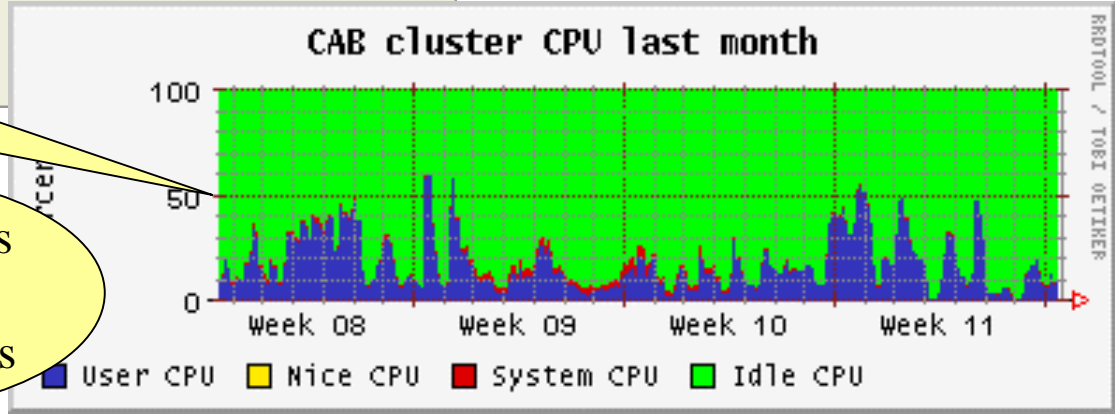
600 CPUs



FNAL-farm Reconstruction Farm

CAB Usage will increase dramatically in the coming months

50% Utilization
Central-Analysis Backend Compute Servers





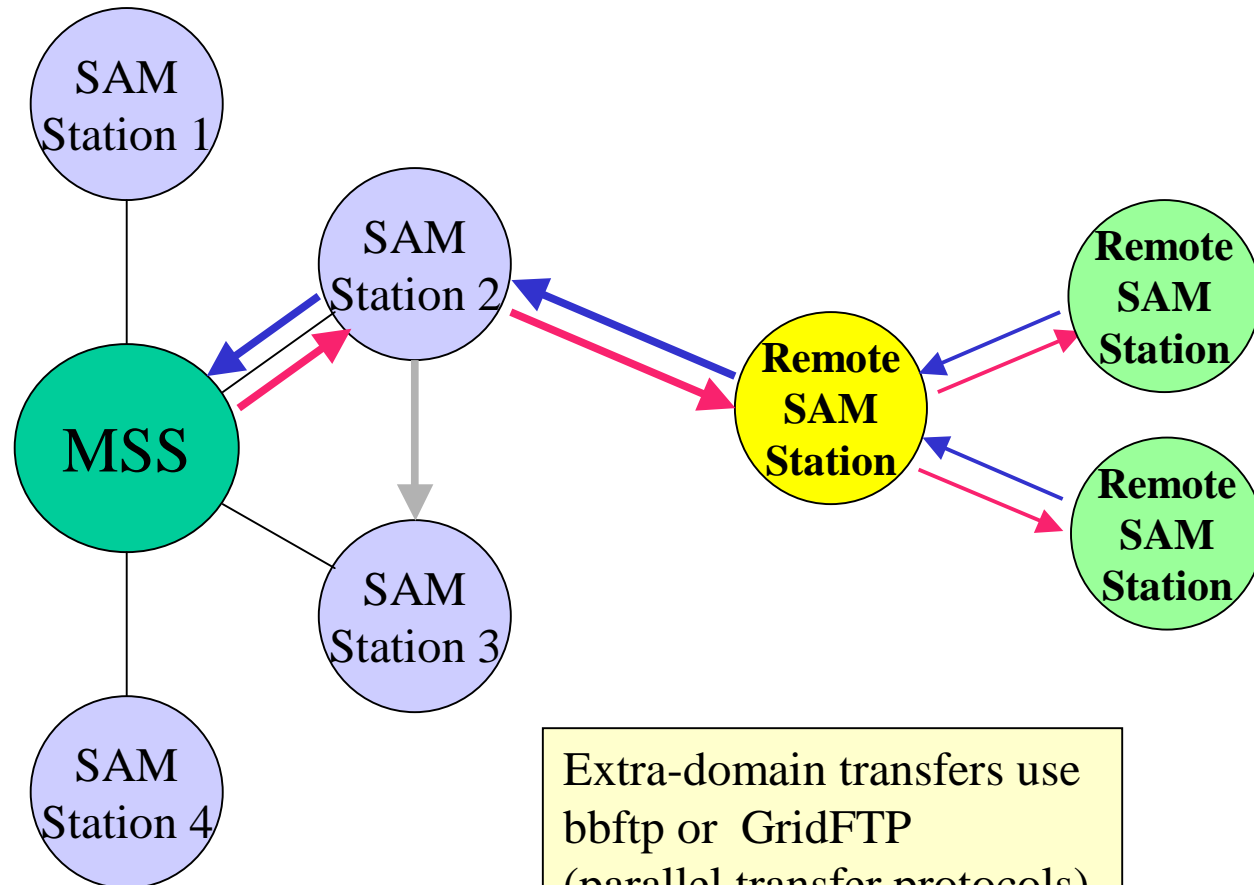
Data to and from Remote Sites

Data Forwarding and Routing



Station Configuration

- **Replica location**
 - Prefer
 - Avoid
- **Forwarding**
 - File stores can be forwarded through other stations
- **Routing**
 - Routes for file transfers are configurable

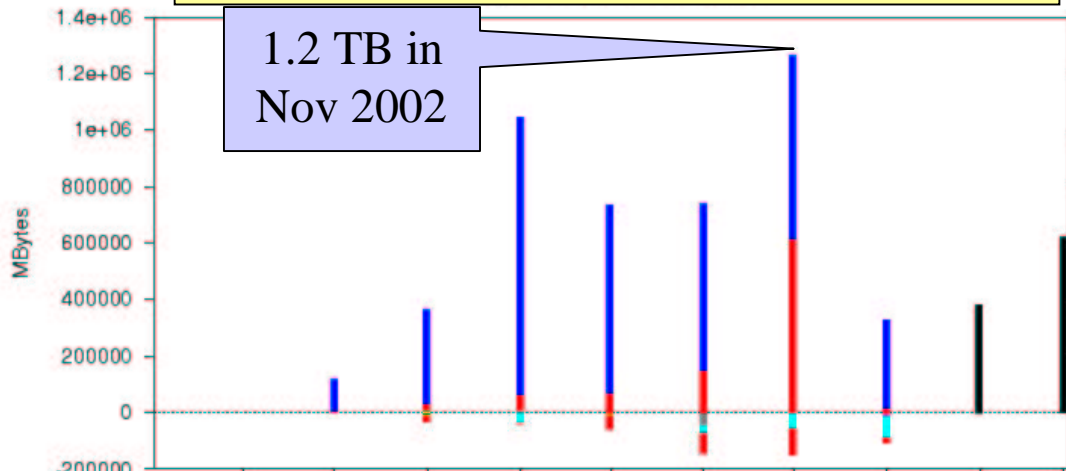




DØ Karlsruhe Station at GridKa



Monthly Thumbnail Data Moved to GridKa

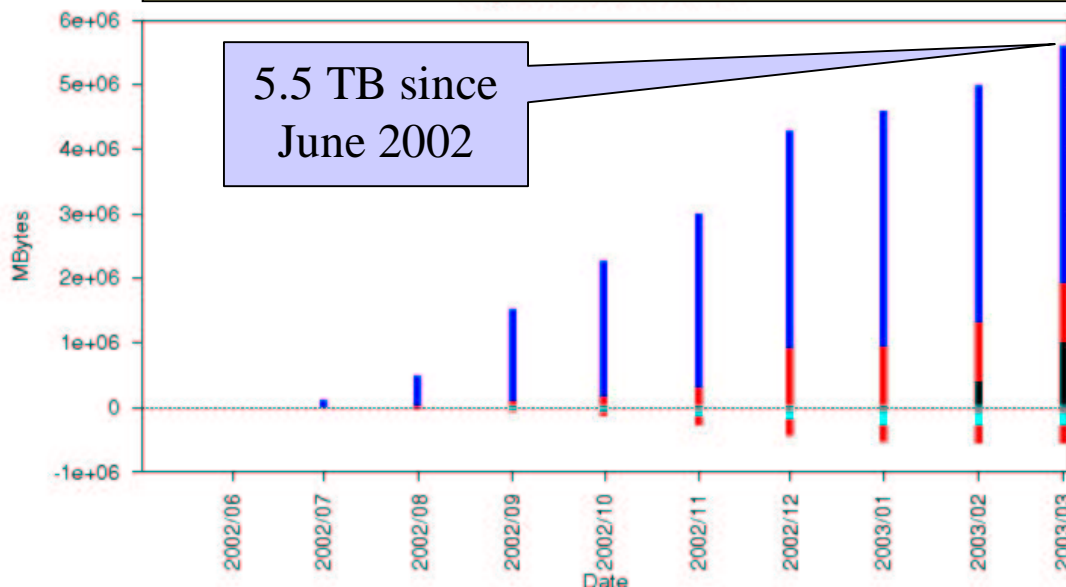


The GridKa SAM Station uses shared cache config. with workers on a private network



This is our first Regional Analysis Center (RAC).

Cumulative Thumbnail Data Moved to GridKa



Resource Overview:

- Compute: 95 x dual PIII 1.2GHz, 68 x dual Xeon 2.2 GHz. DØ requested 6%. (updates in April)
- Storage: DØ has 5.2 TB cache. Use of % of ~100TB MSS. (updates in April)
- Network: 100Mb connection available to users.
- Configuration: SAM w/ shared disk cache, private network, firewall restrictions, OpenPBS, Redhat 7.2, k 2.418, DØ software installed.



Challenges



- Getting SAM to meet the needs of DØ in the many configurations is and has been an enormous challenge. Some examples include...
 - **File corruption issues**. Solved with CRC.
 - **Preemptive distributed caching** is prone to race conditions and log jams. These have been solved.
 - **Private networks** sometimes require “border” naming services. This is understood.
 - **NFS shared cache configuration** provides additional simplicity and generality, at the price of scalability (star configuration). This works.
 - **Global routing** completed.
 - **Installation procedures** for the station servers have been quite complex. They are improving and we plan to soon have “push button” and even “opportunistic deployment” installs.
 - **Lots of details** with opening ports on firewalls, OS configurations, registration of new hardware, and so on.
 - **Username clashing issues**. Moving to GSI and Grid Certificates.
 - **Interoperability with many MSS**.
 - **Network attached files**. Sometimes, the file does not need to move to the user.





SAM Grid

<http://www-d0.fnal.gov/computing/grid/>



May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.

17



DØ Objectives of SAM-Grid

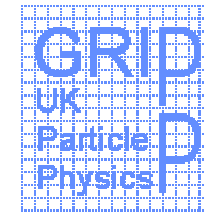


- JIM (Job and Information Management) complements SAM by adding job management and monitoring to data handling.
- Together, JIM + SAM = SAM-Grid
- Bring standard grid technologies (including Globus and Condor) to the Run II experiments.
- Enable globally distributed computing for DØ and CDF.

•People involved:

–Igor Terekhov (FNAL; JIM Team Lead), Gabriele Garzoglio (FNAL), Andrew Baranovski (FNAL), Rod Walker (Imperial College), Parag Mhashilkar & Vijay Murthi (via Contr. w/ UTA CSE), Lee Lueking (FNAL; Team rep. For DØ to PPDG)

–Many others at many DØ and CDF sites



Condor
High Throughput Computing



May 12-15, 2003

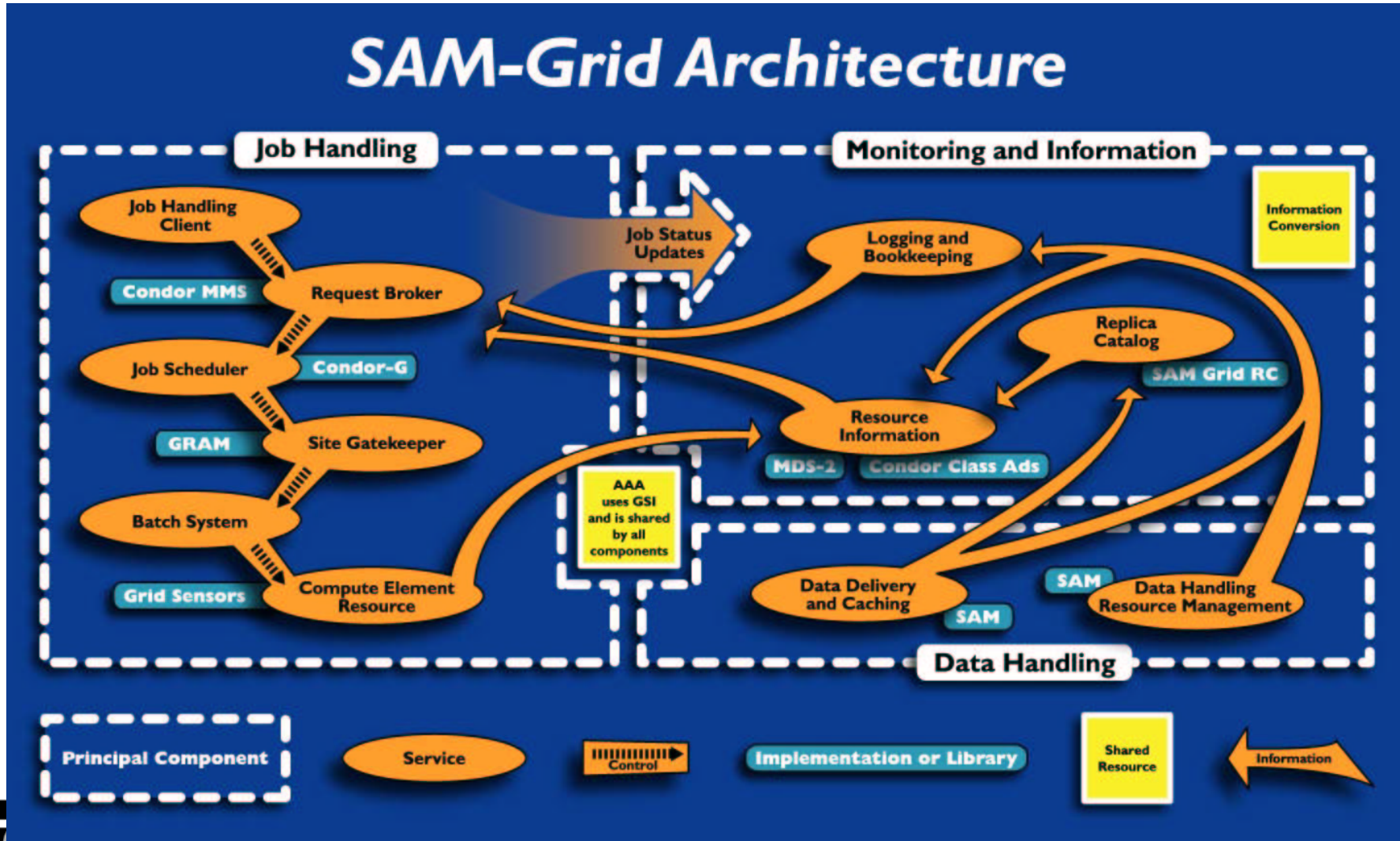
Lee Lueking, EDG Int. Proj. Conf.



the globus project™
www.globus.org



The SAM-Grid Architecture





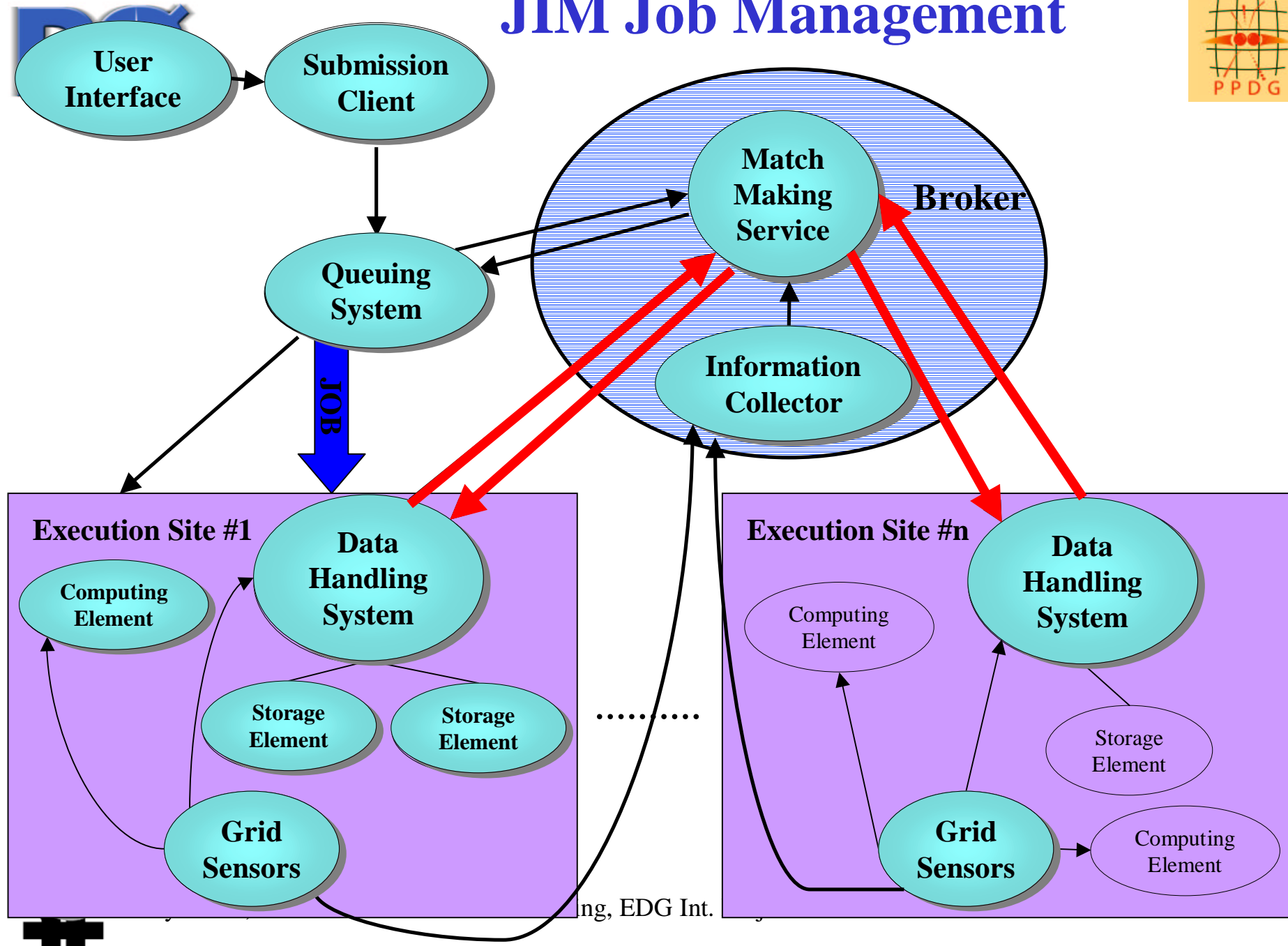
Condor-G Extensions Driven by JIM



- **The JIM Project team has inspired many Extensions to the Condor software**
 - **Added Match Making to the Condor-G for grid use.**
 - **Extended class adds to have the ability to call external functions from the match making service.**
 - **Introduced a three tier architecture which separates the user submission, job management service, and submission sites completely.**
- **Decision making on the grid is very difficult. The new technology allows:**
 - **Including logic not expressible in class ads**
 - **implementing very complex algorithms to establish ranks for the jobs in the scheduler**
- **Also, many robustness and security issues have been addressed**
 - **TCP replaces UDP for communication among Condor services**
 - **GSI now permeates the Condor-G services, driven by the requirements of the three-tier architecture**
 - **Re-matching a grid job that failed during submission**



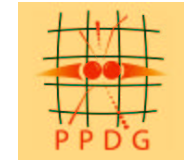
JIM Job Management



ing, EDG Int.



SAM-Grid Monitoring



SAM GRID INFORMATION AND MONITORING SYSTEM - Microsoft Internet Explorer

Address: http://sanadans.fnal.gov:8060/prototype/

SAM GRID INFORMATION & MONITORING SYSTEM

Launching the Monitoring Systems:
Please click at the map to monitor the execution sites.
Click [here](#) to get information about the submission sites.

Participating Experiments:
● D0 ● CDF

SAM Grid Monitoring System - Microsoft Internet Explorer

Address: http://sanadans.fnal.gov:8060/prototype/fnwn_san_stations.php?site=IC

SAM Grid Monitoring System

Wed, 5 Mar 2003 19:56:57 -0600

SAM Grid Monitoring System - Microsoft Internet Explorer

Address: http://sanadans.fnal.gov:8000/prototype/fnwn_san_stations.php?site=FNAL

Monitoring at the FNAL Site

View Authorized Grid Users

Please click on a station's name to get its Server Version and Start-time.
For stations that are grid-enabled, the Cluster Details can be viewed through the available link.

Station Name	Universe	Grid-enabled	Projects	Disks	Groups	Experiment
sanadans	dev	Yes	0	3	9	d0
sammy	dev	Yes	0	2	4	d0
samegg	dev	Yes	0	2	4	d0
control-g	dev	Yes	0	2	4	d0
cdf-glass	dev	Yes	0	2	4	d0
droid	dev	Yes	0	2	4	d0
fnal-farm	dev	Yes	0	2	4	d0
cdf-glass	dev	Yes	0	2	4	d0
cdf-glass	dev	Yes	0	2	4	d0

SAM Grid Projects - Microsoft Internet Explorer

Projects at: fnal-farm - prd

Sam Project Id	Total Files	Locked	Given	Delivery Errors	Wanted	Local Owner	Group
farm.r13.06.01.23325	0	0	0	0	0	d0farm	d0production
farm.p13.06.01.23345	0	0	0	0	0	d0farm	d0production
farm.r13.06.01.23347	0	0	0	0	0	d0farm	d0production
farm.r13.06.01.23348	0	0	0	0	0	d0farm	d0production
farm.r13.06.01.23351	0	0	0	0	0	d0farm	d0production
farm.r13.06.01.23354	5	0	5	0	0	d0farm	d0production
farm.r13.06.01.23355	0	0	0	0	0	d0farm	d0production
farm.r13.06.01.23356	2	0	2	0	0	d0farm	d0production
farm.r13.06.01.23357	0	0	0	0	0	d0farm	d0production
farm.r13.06.01.23358	26	0	26	0	0	d0farm	d0production
farm.r13.06.01.23359	29	0	29	0	0	d0farm	d0production
farm.r13.06.01.23360	2	0	2	0	0	d0farm	d0production
farm.r13.06.01.23363	12	0	12	0	0	d0farm	d0production
farm.r13.06.01.23362	78	0	78	0	0	d0farm	d0production
farm.r13.06.01.23364	100	1	80	0	20	d0farm	d0production
farm.r13.06.01.23365	99	0	78	0	21	d0farm	d0production
farm.r13.06.01.23366	56	0	56	0	0	d0farm	d0production
farm.r13.06.01.23367	99	0	45	0	54	d0farm	d0production
farm.r13.06.01.23368	56	0	20	0	36	d0farm	d0production
farm.r13.06.01.23369	99	0	20	0	79	d0farm	d0production
farm.r13.06.01.23370	57	11	21	0	36	d0farm	d0production
farm.r13.06.01.23374	5	0	5	0	0	d0farm	d0production
farm.r13.06.01.23376	7	0	7	0	0	d0farm	d0production
farm.r13.06.01.23375	2	0	2	0	0	d0farm	d0production
farm.r13.06.01.23377	2	0	2	0	0	d0farm	d0production
farm.r13.06.01.23386	11	0	11	0	0	d0farm	d0production



May 12-15, 2003

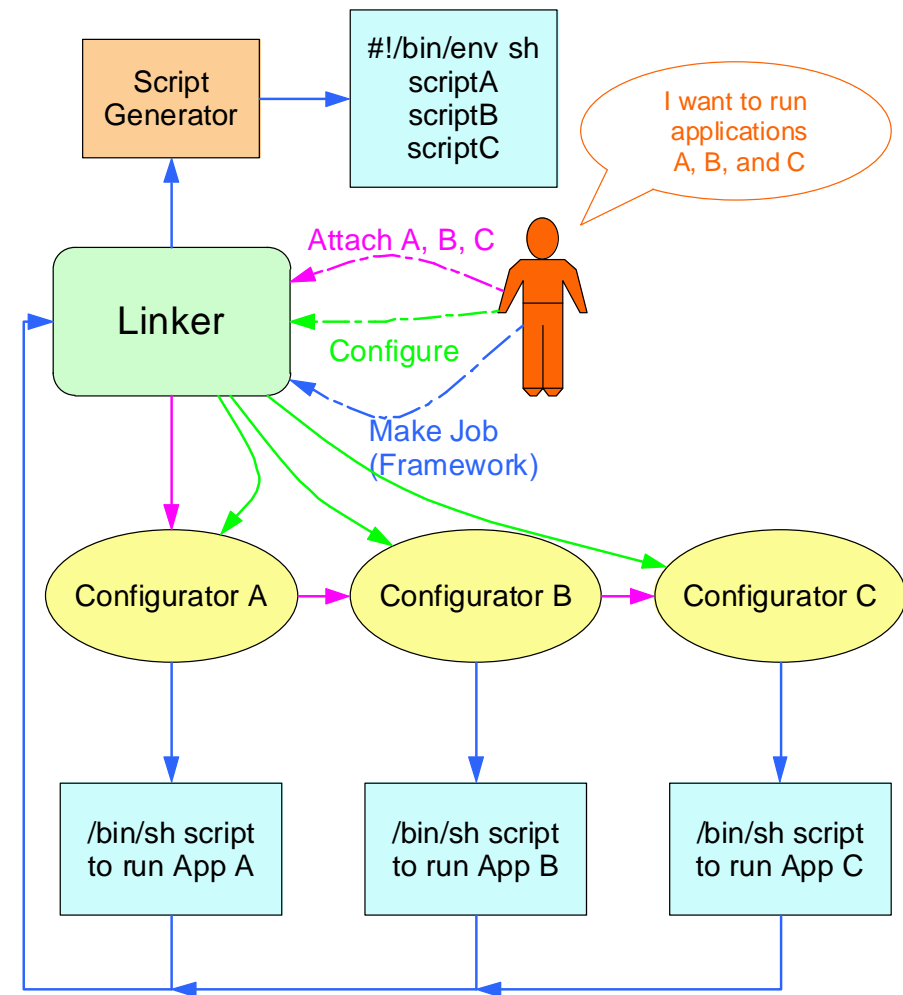
Lee Lueking, EDG Int. Proj. Conf.



Meta Systems



- MCRunJob approach by CMS and DØ production teams
- Framework for dealing with multiple grid resources and testbeds (EDG, IGT)





DO JIM Deployment



- **A site can join SAM-Grid with combinations of services:**
 - **Monitoring, and/or**
 - **Execution, and/or**
 - **Submission**
- **May 2003: Expect 5 initial execution sites for SAMGrid deployment, and 20 submission sites.**
- **Summer 2003: Continue to add execution and submission sites.**
- **Grow to dozens execution and hundreds of submission sites over next year(s).**
- **Use grid middleware for job submission within a site too!**
 - **Administrators will have general ways of managing resources.**
 - **Users will use common tools for submitting and monitoring jobs everywhere.**





What's Next for SAM-Grid?



After JIM version 1

- **Improve scheduling jobs and decision making.**
- **Improved monitoring, more comprehensive, easier to navigate.**
- **Execution of structured jobs**
- **Simplifying packaging and deployment. Extend the configuration and advertising features of the uniform framework built for JIM that employs XML.**
- **CDF is adopting SAM and SAM-Grid for their Data Handling and Job Submission.**
- **Co-existence and Interoperability with other Grids**
 - **Moving to Web services, Globus V3, and all the good things OGSA will provide. In particular, interoperability by expressing SAM and JIM as a collection of services, and mixing and matching with other Grids**
 - **Work with EDG and LCG to move in common directions**





Run II plans to use the Virtual Data Toolkit



- JIM is using advanced version of Condor-G/Condor - actually driving the requirements. Capabilities available in VDT 1.1.8 and beyond.
- D0 uses very few VDT packages- Globus GSI, GridFTP, MDS and Condor.
- JIM ups/upd packaging includes configuration information to save local site managers effort. Distribution and configuration tailored for existing/long legacy D0 systems.
- Plans to work with VDT such that D0-JIM will use VDT in the next six months.
- ==>> VDT versions are currently being tailored for each application community. This cannot continue. We - D0, US CMS, PPDG, FNAL, etc.- will work with the VDT team and the LCG to define how VDT versions should be
 - Constructed and Versioned
 - Configured
 - Distributed to the various application communities
 - Requirements and scheduled for releases.

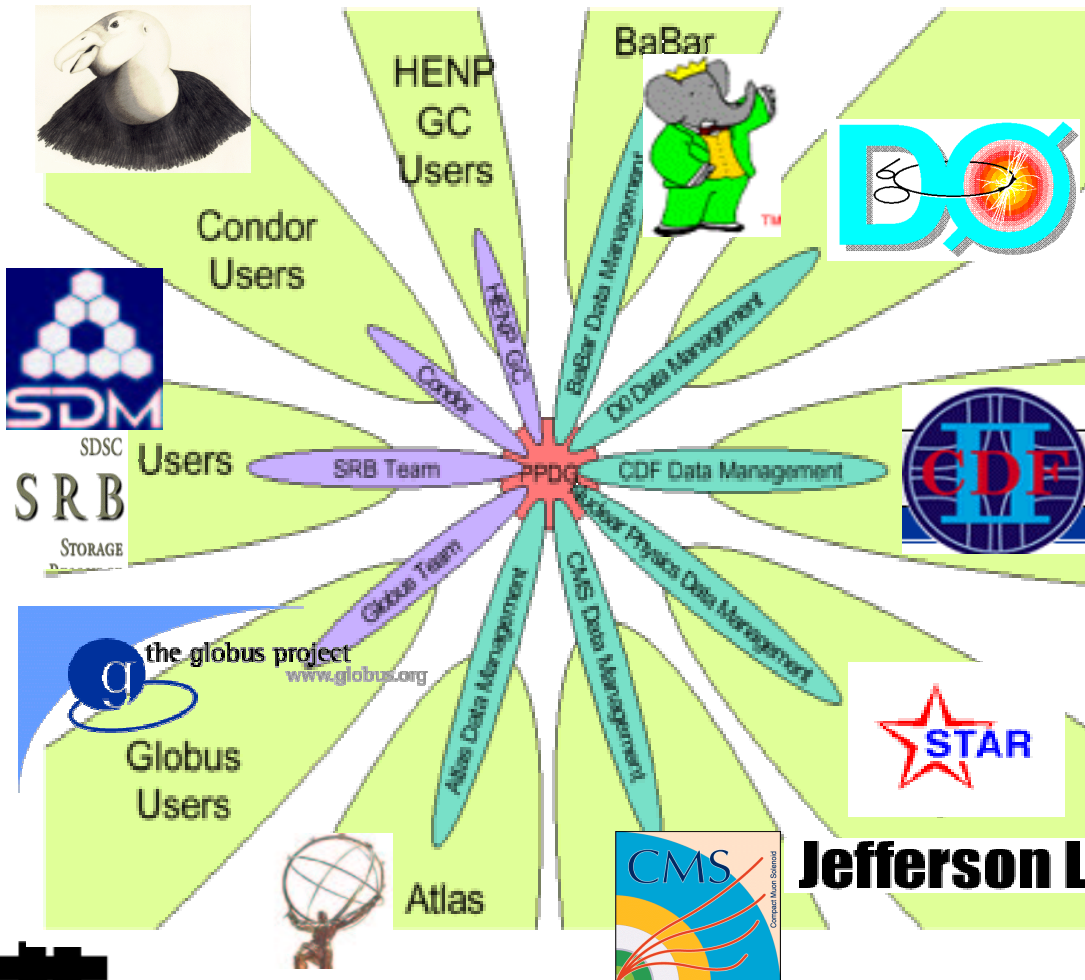




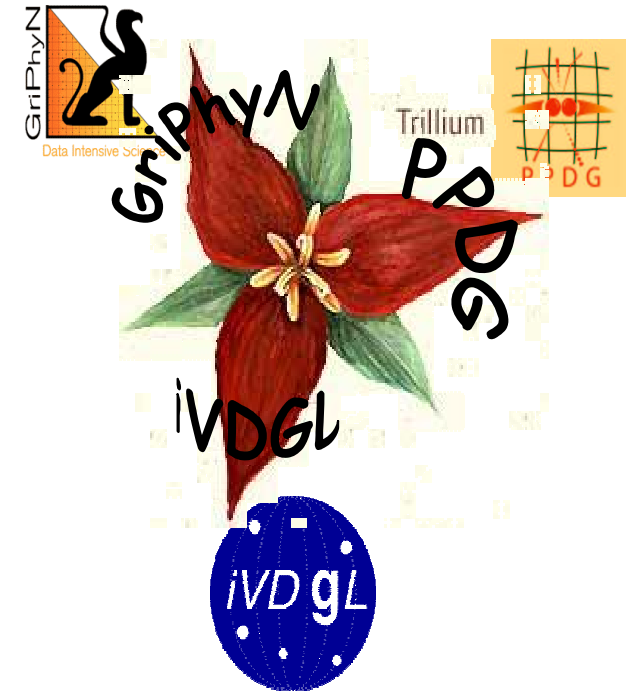
Projects Rich in Collaboration



PPDG



Trillium



May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.



Collaboration between Run 2 and US CMS Computing at Fermilab



- D0, CDF, and CMS are all using Dcache and Enstore storage management systems.
- Grid VO management - joint US-CMS, iVDGL, INFN-VOMS, (LCG?) project is underway
 - <http://www.uscms.org/s&c/VO/meeting/meet.html>
 - There is a commitment from the RUN II Experiments to collaborate on with this effort in near future.
- (mc)Runjob scripts - joint work on core framework between CMS and Run II experiments has been proposed.
- Distributed and Grid accessible databases and applications are a common need.
- As part of PPDG we expect to collaborate on future projects such as Troubleshooting Pilots (end to end error handling and diagnosis).
- Common infrastructure in Computing Division for system and core service support etc. ties us together.





Regional Computing Approach



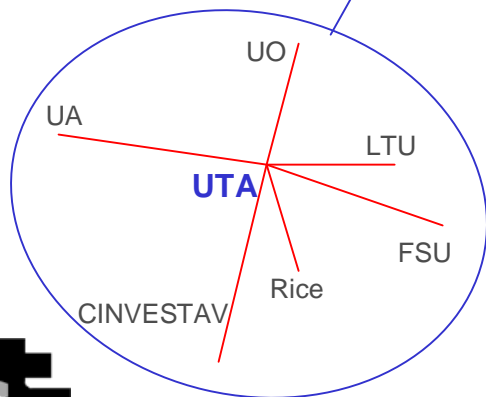
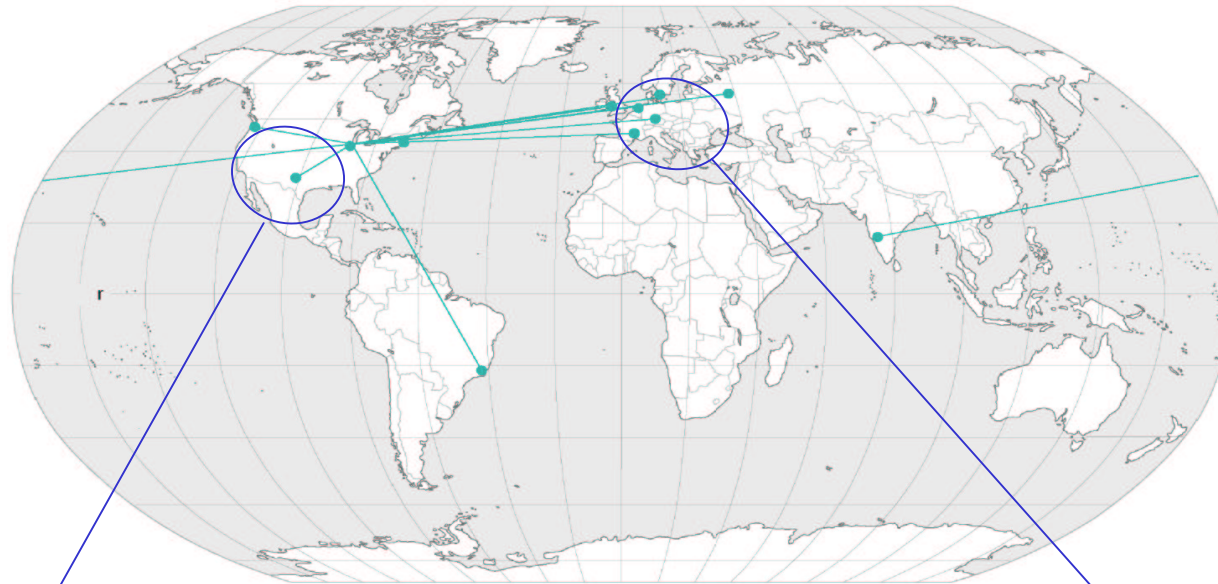
May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.

29



DØ Regional Model



Centers also in the UK and France

UK: Lancaster, Manchester, Imperial College, RAL

France: CCin2p3, CEA-Saclay, CPPM Marseille, IPNL-Lyon, IRES-Strasbourg, ISN-Grenoble, LAL-Orsay, LPNHE-Paris



May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.



Regional Analysis Centers (RAC) Functionality



- **Preemptive caching**
 - **Coordinated globally**
 - **All DSTs on disk at the sum of all RAC's**
 - **All TMB files on disk at all RACs, to support mining needs of the region**
 - **Coordinated regionally**
 - **Other formats on disk: Derived formats & Monte Carlo data**
- **On-demand SAM cache: ~10% of total disk cache**
- **Archival storage (tape - for now)**
 - **Selected MC samples**
 - **Secondary Data as needed**
- **CPU capability**
 - **supporting analysis, first in its own region**
 - **For re-reconstruction**
 - **MC production**
 - **General purpose DØ analysis needs**
- **Network to support intra-regional, FNAL-region, and inter-RAC connectivity**

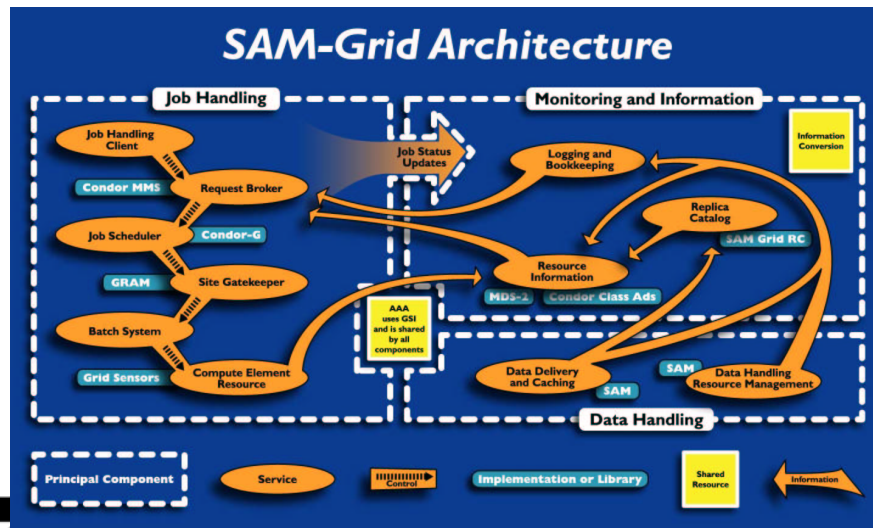
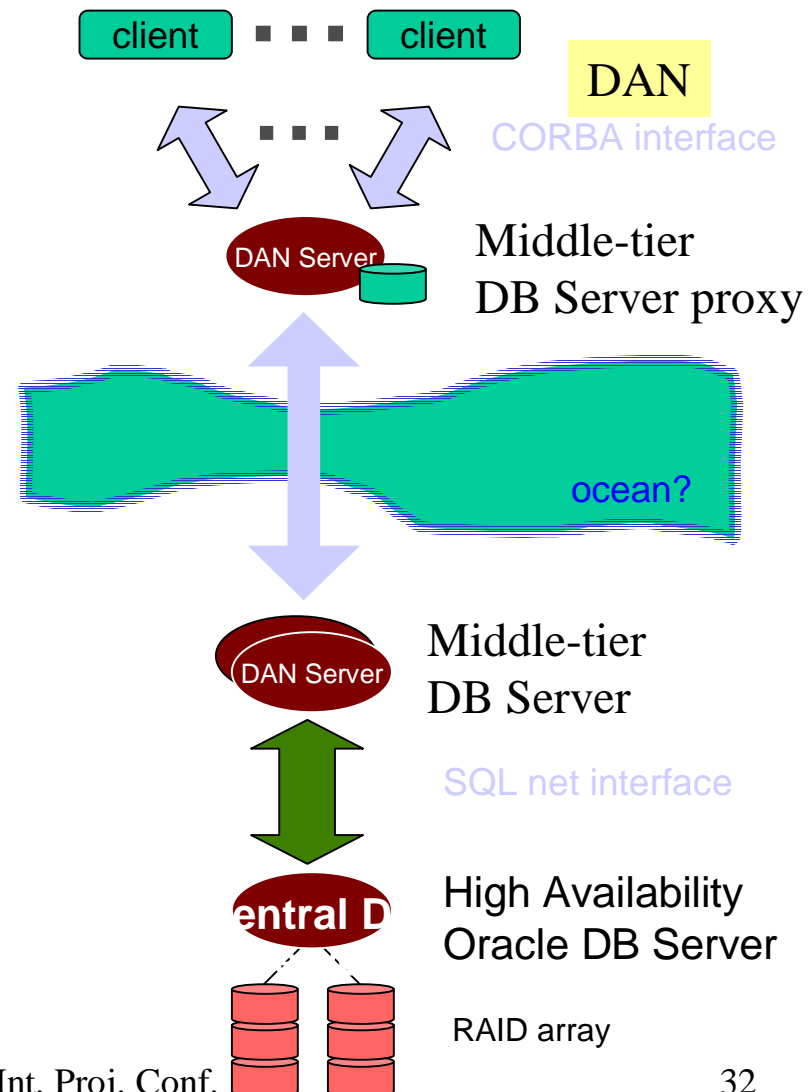




Required RAC Server Infrastructure



- SAM-Grid Gateway machine
- Oracle database access servers
 - Provided via middle tier server (DAN)
 - DAN = Database Access Network
- Accommodate realities like:
 - Policies and culture for each center
 - Sharing with other organizations
 - Firewalls, private networks, et cetera





Summary of Current & Soon-to-be RACs



Regional Centers	Institutions within Region	CPU ΣHz (Total*)	Disk (Total*)	Archive (Total*)	Schedule
GridKa @FZK	Aachen, Bonn, Freiburg, Mainz, Munich, Wuppertal,	52 GHz (518 GHz)	<div style="border: 1px solid black; background-color: #f4a460; padding: 10px; text-align: center;"> Total Remote CPU 360 GHz (1850 GHz) </div>		Established as RAC
SAR @UTA (Southern US)	AZ, Cinvestav (Mexico City), LA Tech, Oklahoma, Rice, KU, KSU	160 GHz (320 GHz)			Summer 2003
UK @tbd	Lancaster, Manchester, Imperial College, RAL	46 GHz (556 GHz)			Active, MC production
IN2P3 @Lyon	CCin2p3, CEA-Saclay, CPPM-Marseille, IPNL-Lyon, IRES-Strasbourg, ISN-Grenoble, LAL-Orsay, LPNHE-Paris	100 GHz			Active, MC production
DØ @FNAL (Northern US)	Farm, cab, clued0, Central-analysis	1800 GHz	<div style="border: 1px solid black; background-color: #00ffff; padding: 10px; text-align: center;"> FNAL CPU 1800 GHz </div>		Established as CAC

*Numbers in () represent totals for the center or region, other numbers are DØ's current allocation.





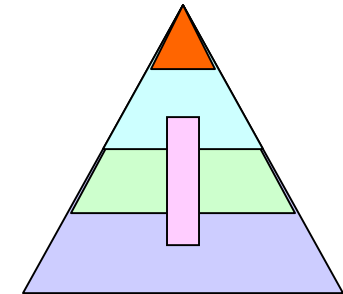
Data Model



Fraction of Data Stored

Data Tier	Size/event (kB)	FNAL Tape	FNAL Disk	Remote Tape	Remote Disk
RAW	250	1	0.1	0	0
Reconstructed	50	0.1	0.01	0.001	0.005
DST	15	1	0.1	0.1	0.1
Thumbnail	10	4	1	1	2
Derived Data	10	4	1	1	1
MC D0Gstar	700	0	0	0	0
MC D0Sim	300	0	0	0	0
MC DST	40	1	0.025	0.025	0.05
MC TMB	20	1	1	0	0.1
MC PMCS	20	1	1	0	0.1
MC root-tuple	20	1	0	0.1	0
Totals RIIa ('01-'04)/ RIIb ('05-'08)		1.5PB/ 8 PB	60TB/ 800 TB	~50TB	~50TB

per Region
Data Tier
Hierarchy



▲ Metadata
~0.5TB/year

Numbers are
rough estimates

the cpb model presumes:
25Hz rate to tape, Run IIa
50Hz rate to tape, Run IIb
events 25% larger, Run IIb





Challenges



- Operation and Support
 - Ongoing shift support: 24/7 “helpdesk” shifters (trained physicists)
 - SAM-Grid station administrators: Expertise based on experience installing and maintaining the system
 - Grid Technical Team: Experts in SAM-Grid, DØ software + technical experts from each RAC.
 - Hardware and system support provided by centers
- Production certification
 - All DØ MC, reconstruction, and analysis code releases have to be certified
- Special requirements for certain RAC’s
 - Forces customization of infrastructure
 - Introduces deployment delays
- Security issues, grid certificates, firewalls, site policies.





Operations



Expectation Management





Summary



- The DØ Experiment is moving toward exciting Physics results in the coming years.
- The software is stable and provides reliable data delivery and management to production systems worldwide.
- SAM-Grid is using standard Grid middleware to enable complete Grid functionality. This is rich in collaboration with Computer Scientists and other Grid efforts.
- DØ will rely heavily on remote computing resources to accomplish its Physics goals





Thank You



May 12-15, 2003

Lee Lueking, EDG Int. Proj. Conf.

38