

DO Computing Experience

Amber Boehnlein

FNAL/CD

For DO collaborations

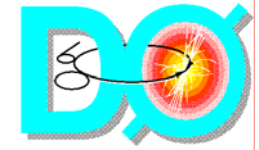
April 28, 2005

Introduction



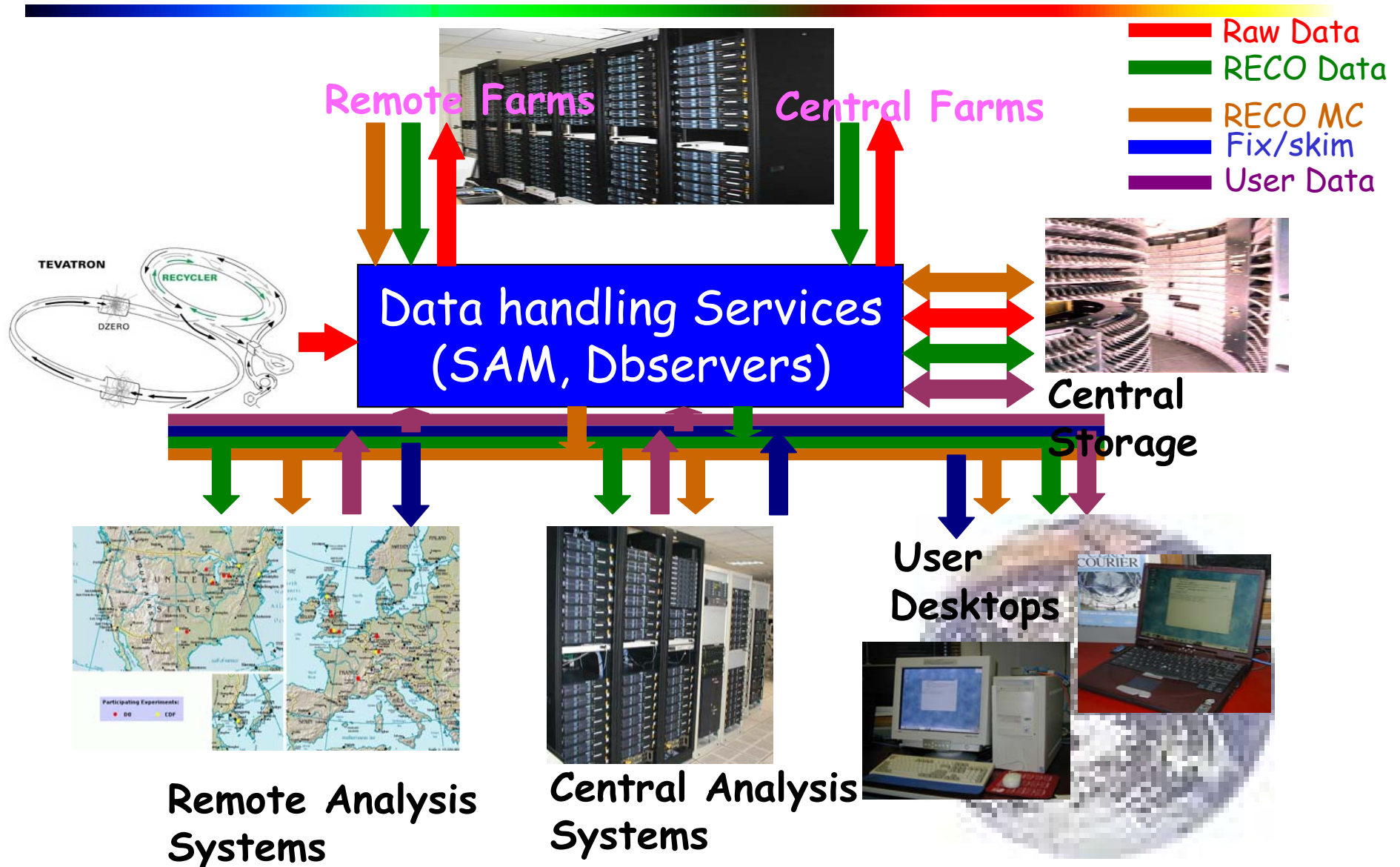
- My understanding of the charge of this series of conferences is to understand how Tevatron experience can be leveraged for the LHC experiments.
 - ◆ Congratulations to the organizers.
- For LHC computing and software, most of the technical choices are made and any re-evaluations should be made in the context of the LHC computing experience
 - ◆ This is an exciting time for LHC computing, the challenges are already apparent
 - ◆ Many Run II computing experts working on LHC computing
- Focus on how the 1997 plan evolved into 2005 reality
 - ◆ I led/co-led DO Computing and Software from 2001-2004
 - ◆ Currently my FNAL CD role is to run what one could consider the RunII Tier 0 center
 - ◆ Try to avoid many detailed examples- case studies
- DO computing is extremely successful
 - ◆ Credit for that belongs to the 1997 pioneers: Wyatt Merritt, Vicky White, Lee Lueking, Heidi Schellman, Mike Diesburg, Stu Fuess
 - ◆ Credit also belongs to the second generation, particularly the European community who committed to the global computing model with enthusiasm
 - ◆ Credit also to my successor in the DO experiment role, Gustaaf Brooijmans

Historical Perspective



- **Planning for Run II computing started in approximately 1996.**
- **Both experiments took a critical look at the Run I computing and software models. Run I computing largely met the needs of the experiments, but wouldn't scale to Run II rates and had known deficiencies**
- **Input from physics groups**
- **The result was a bottoms up need estimate that drove the specification of experiment specific algorithm software, supporting infrastructure and data handling and interaction with storage elements.**
- **Job management and workflow not part of initial thinking**
- **The basic elements on the computing side are as envisioned by the plan.**
- **I'm not going to talk much about the offline**

Computing Model

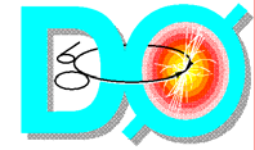


Vital Statistics



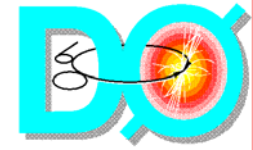
DO Vital Statistics	1997	2005
Peak (Average) Data Rate(Hz)	50(20)	50(17)
Events Collected	600M/year	1.2 B
Raw Data Size (kbytes/event)	250	250
Reconstructed Data Size (kbytes/event)	100 (5)	200 (20->60)
User formats (kbytes/event)	1	40
Tape storage	280 TB/year	1 pb on tape
Tape Reads/writes (weekly)		30TB/7TB
Analysis/cache disk	7TB/year	120 TB
Reconstruction Time (Ghz-sec/event)	2.00	50 (120)
Monte Carlo Chain	full Geant	full Geant
user analysis times (Ghz-sec/event)	?	1
user analysis weekly reads	?	1B events
Primary Reconstruction farm size (GHz)	50	550000
Central Analysis farm size (GHz)	50	500000
Remote resources(GHz)	?	2500000

The Program of Work



- In order to realize the 1997 plans, a number of joint projects and working groups were undertaken by CDF, DO and FNAL CD and other parties
 - ◆ SoftRelTools
 - ◆ C++ support and libraries
 - ◆ Sequential Access Via Metadata (SAM)*
 - ◆ Enstore (storage interface)
 - ◆ Many others
 - ◆ JIM (SAMGrid)
 - ◆ dCache*
- Collaborations on ROOT
- CD assignments for experiment specific projects, framework, event data model, databases, online
- Experimenters assignments for detector software, MC, algorithms, object id
- About 30 FTEs from FNAL CD—plus 60 FTE (my guess) experiment for offline development at peak
 - ◆ Most efforts under manned and under directed
- Currently, 15 FTE direct support from FNAL CD + 30 FTE (self- effort report) for Computing support today
 - ◆ Note: “Misc. management” self-effort estimate comparable
- In Retrospect
 - ◆ The basis for DO’s computing is founded on the Run II working groups and joint projects
 - ◆ I personally think that the joint projects have been very successful, but the “joint” aspect has not always been realizable.

SAM Data Handling

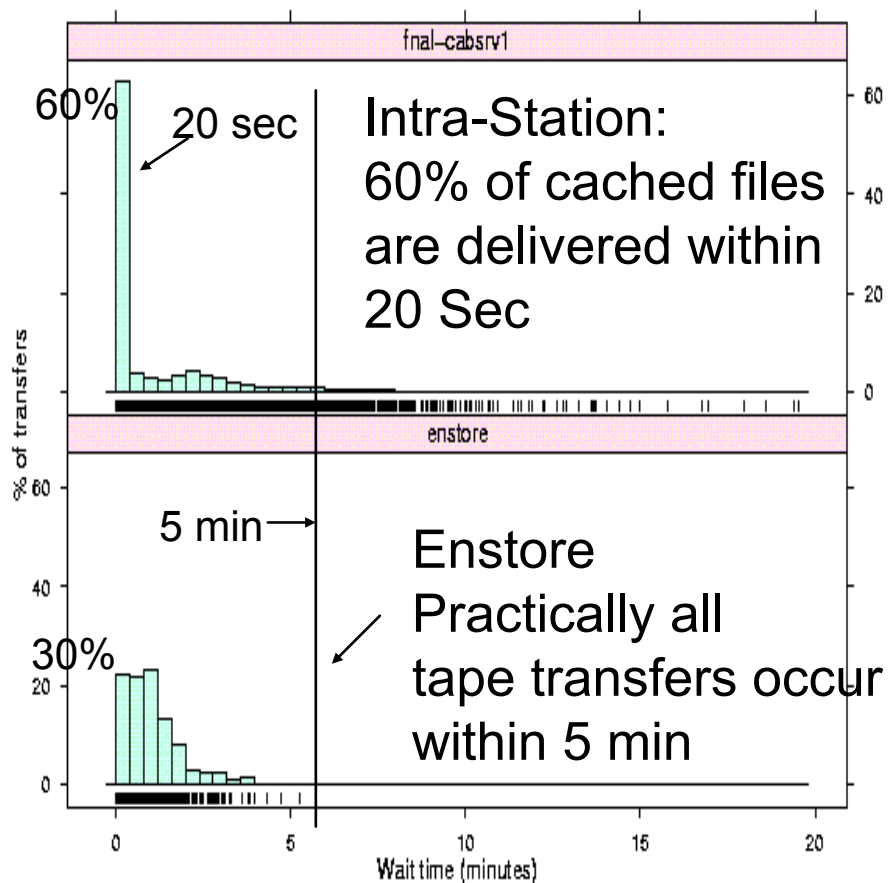


- **Flagship CD-Tevatron Joint project—initial design work ~7 years ago, in production for DO for 5 years**
- **Provides transparent global access to the data**
- **Stable SAM operations allows for global support and additional development**
- **Services provided**
 - ◆ **Comprehensive meta-data to describe collider and Monte Carlo data.**
 - ◆ **Bookkeeping services**
 - ◆ **Consistent user interface via command line and web**
 - ◆ **Local and wide area data transport**
 - ◆ **Caching layer**
 - ◆ **Batch adapter support (PBS, Condor, Isf, site-specific batch systems)**
 - ◆ **Optimization knobs in place**
- **Second Generation –Experience and new perspectives extend and improve functionality**
 - ◆ **Schema and DBserver updated in 2004**
 - ◆ **Introduction of SRM interface/dCache**
 - ◆ **Monitoring and Information Server prototype**
 - **move away from log file monitoring**
 - **Provide more real time monitoring**

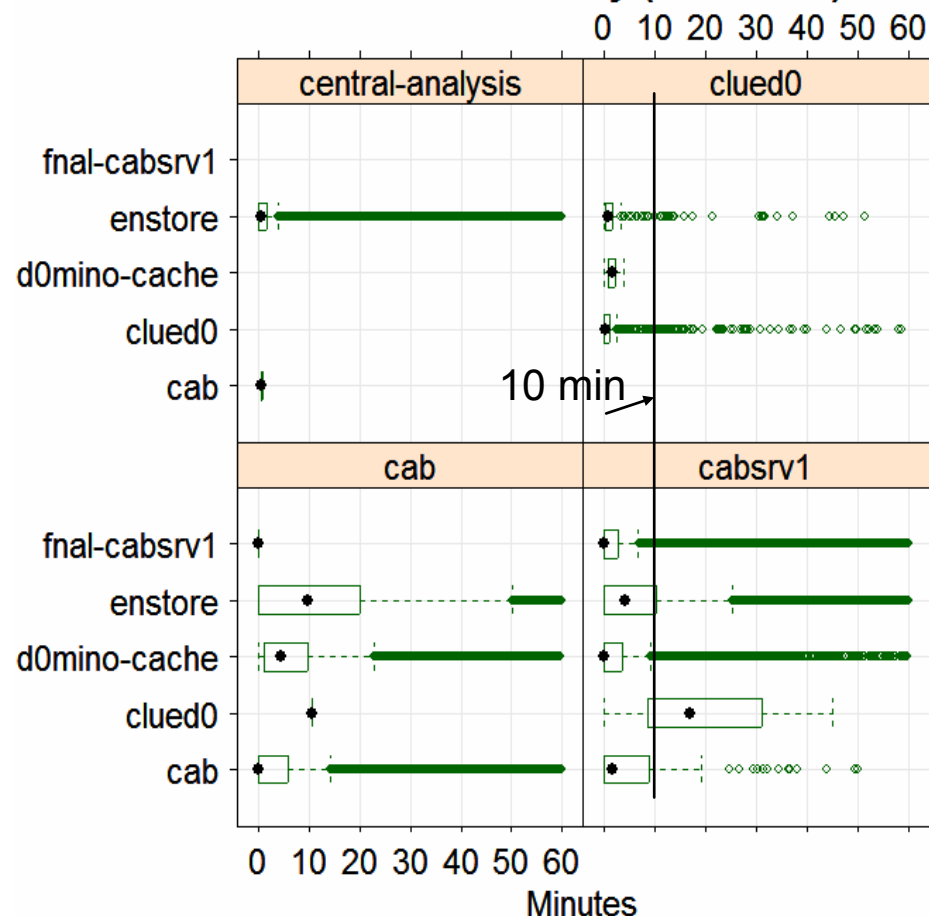
DO SAM Performance



Process Wait Times

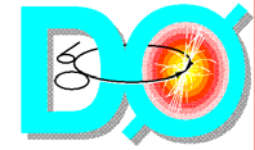


Wait Time for File Delivery (truncated)



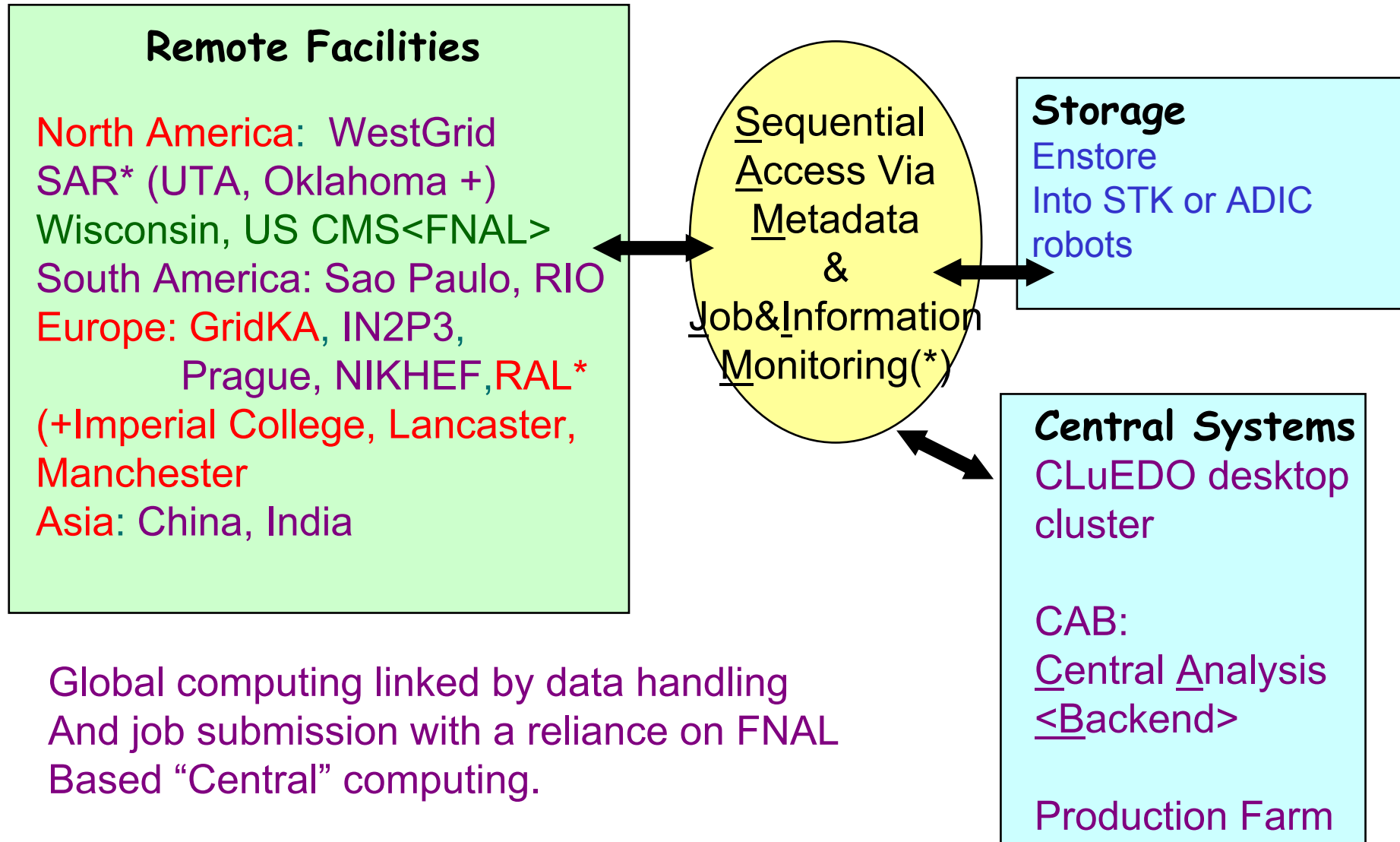
Before adding 100 TB of Cache, 2/3 transfers could be from tape.
Still robust!

SAMGrid



- **SAMGrid project (started around 2001) includes Job and Information Monitoring (JIM), grid job submission and execution package**
 - ◆ **Uses VDT for underlying services**
 - ◆ **JIM is in production for execution at several DO MC sites**
 - ◆ **JIM is being used for reprocessing for 5 certified sites, 3-4 in progress**
 - ◆ **Collaboration/discussions within the experiments on the interplay of LCG and Open Science Grid with SAMGrid efforts**
 - **Demonstration of use of sam_client on LCG site**
 - **University of Oklahoma runs Grid3 and JIM on a single gatekeeper (co-existence...)**
 - **Reprocessing and MC generation at Wisconsin, CMS resources in co-existence mode**
 - **Interoperability will be a focus this summer**
 - **Will be working closely with LHC people**

DO Computing today



Global computing linked by data handling
And job submission with a reliance on FNAL
Based “Central” computing.

Oversight by Computing Planning Board

Amber Boehnlein, FNAL

Virtual Center



- In order to assign a money value for remote computing contributions to DO, we developed a model based on a virtual center
- The Virtual Center represents the system that would be needed to supply ALL DO computing needs
 - ◆ Cost if all computing is at FNAL
 - ◆ purchased in the yearly currency
 - Disk and servers and CPU for FNAL analysis
 - Production activities such as MC generation, processing and reprocessing.
 - Mass storage, cache machines and drives to support extensive data export
- Assign fractional value for remote contributions based on fractions of events processed * the value of the entire task
 - ◆ Assigning equipment purchase cost as value (“Babar Model”) doesn’t take into account life cycle of equipment nor system efficiency
 - ◆ DO Computing planning board includes strong remote participation, representation—necessary to schedule and secure resources and interface to worldwide computing centers.
 - ◆ In general, the resources are not specifically assigned to DO.
- Value estimates are work in progress
 - ◆ No value assigned Wide Area Networking, Infrastructure, desktop computing, analysis
 - ◆ Requesting resources and keeping to schedules are challenging.

Budget Request



- **Initial Run II funding for computing delayed**
 - ◆ Initial estimates eye-popping relied on SMP machines.
 - ◆ Delayed funding was ok—Moore’s law, cheap processors and file servers and schedule slip covered shortfalls of estimates.
 - ◆ Budgets falling again, making hard choices.
- **Yearly bottoms-up estimate of equipment budget**
 - ◆ Some top-down budget guidance
 - ◆ Estimates use past year’s experience
 - ◆ Externally reviewed.
- **Budgets falling for the past few years**
 - ◆ Fitting \$1.8M worth of computing requests into \$1.25M actual budget

	Purchased 2003	Purchased 2004	Purchase 2005	Purchase 2006	Purchase 2007	Purchase 2008
FNAL Analysis CPU	\$470,000	\$277,000	\$417,132	\$534,926	\$406,376	\$350,311
FNAL Reconstruction	\$200,000	\$370,000	\$454,269	\$717,742	\$443,490	\$362,546
File Servers/disk	\$111,000	\$350,000	\$357,000	\$356,000	\$293,000	\$276,000
Mass Storage	\$280,000	\$254,700	\$40,000	\$600,000	\$300,000	\$100,000
Infrastructure	\$244,000	\$140,000	\$547,000	\$200,000	\$200,000	\$200,000
FNAL Total	\$1,305,000	\$1,391,700	\$1,815,402	\$2,408,667	\$1,642,867	\$1,288,856

AMBER DUEHLIN, FNAL

Budget History



	2003			2004	
	Projected	Projected(2)	Purchased	Projected	Purchased
FNAL Analysis CPU	\$505,400	\$500,000	\$470,400	\$339,000	\$277,000
FNAL Reconstruction	\$200,000	\$40,000	\$200,000	\$83,000	\$370,000
File Servers/disk	\$262,000	\$200,000	\$111,000	\$490,000	\$350,000
Mass Storage	\$460,000	\$285,000	\$280,000	\$230,000	\$254,700
Infrastructure	\$640,000	\$500,000	\$244,000	\$290,000	\$140,000
FNAL Total	\$2,067,400	\$1,525,000	\$1,305,400	\$1,432,000	\$1,391,700

Reconstruction costs underestimated-delayed deployment of adequate disk

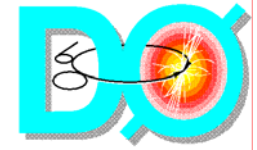
Planning enables an understanding of the trade-offs.

A data handling system that enables use of seamless offsite resources AND prestages data from tape AND robotic storage that out-performs expectation AND network capacity has enabled current budgets to provide sufficient computing for DO

In retrospect: Don't panic-Something will go wrong, but lots of things will go right

Amber Boehlein, FNAL

Case Studies



- Present two examples
- Both examples demonstrate how the collaboration interacted with a set of computing and software decisions and priorities

Case Study: Data Formats

- 1997 plan based on Run I experience, tiered structure of formats
 - ◆ STA raw+ reco for subsample (never implemented)
 - ◆ DST 150Kbytes/event
 - ◆ TMB 10 Kbytes/event
 - ◆ Users expected to produce and maintain analysis specific tuples
 - ◆ Proponents for common tuple format for analysis
- 2001
 - ◆ DST in place, but never caught on with users who preferred to pick raw events
 - ◆ debugging root tuple format produced on farm (RecoCert)
 - Too large and too slow to produce for all events
- Late 2002 introduced the TMB (20 kbytes—included calorimeter energy)
 - TMB extremely successful for data access, algorithm development, but unpacking is slow
 - Slow link times
 - Some users uncomfortable with DO software
 - Some users had significant code base from the debugging tuple
 - ◆ TMBTree introduced as “common” tuple format
 - Some users had significant code base from debugging tuple
 - Backwards compatibility not trivial
 - Did not develop a support base—users view it as a service, while developers view it more as open source.
- All formats
 - ◆ Poorly documented
 - ◆ Poorly thought out
 - ◆ Difficult to maintain and consequently poorly maintained
 - ◆ Spotty support for key applications such as trigger simulation

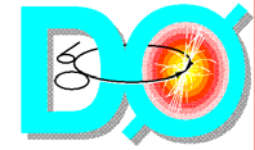
Case Study: Data formats



- 2003 attempted to introduce mechanism to support “Chunks” in root—successfully prototyped-but not pursued.
 - ◆ Late is not always better than than never...
- 2005
 - ◆ Eliminated the DST&TMB—expanded TMB+ with tracking hits fills role in the model
 - ◆ Introducing ROOT “Common Analysis Format” which might or might not fill TMB role
- Striking case of user “stickiness”
 - ◆ During commissioning, many people wanted the shortest learning curve possible--they had work to do NOW.
 - Limited experience with C++
 - No interaction with the release system or code repository
 - No interaction with the D0 framework or code base
 - No interaction with the event data model “Chunks”
 - RecoCert ntuple was the natural choice for many.
 - Created something of a divide between the framework/non-framework users which became SAM/non-SAM users
 - ◆ By 2003, physics groups had distinct patterns end level analysis-complicating standard object ID.
 - ◆ Many have firm, deeply held convictions on the “right” way to handle the data format issue--their format
 - Most people very quickly realize that they don’t want to be a service provider for a global analysis format—the “right” way sometimes includes insisting that “Computing” provide support.

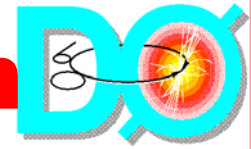
Amber Boehnlein, FNAL

Case Study: Data Formats



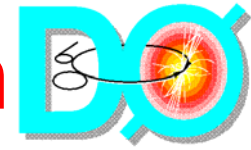
- In retrospect:
 - ◆ In general, the 1997 planning discussions represented the trade-offs quite well, but there still was no consensus on the common analysis format issue.
 - When the official plan wasn't fully realized, the door was open
 - ◆ Introducing the TMB as a well thought out, well organized, well documented data format prior to the start of the run would have been an extremely good use of id and algorithm developer time
 - ◆ Some of the eventual problems could have been sidestepped by an insistence on use of the DST for algorithm development. This would have introduced other problems...
 - ◆ User support starts with design—address infrastructure issues that are barriers.
 - ◆ Schemes that require signing up other people for work they didn't want to do heightened the incentive for “discussion” rather than resolution.
 - People cannot be expected to contribute to tasks that feel (for good or ill) belong to someone else.
 - If we don't build it, they won't come.
 - If we do build it, they still might not come

Case Study: Global Production



- Early SAM design motivated by FNAL based data handling
- Realized early that remote capability was a relatively straightforward way to add value
- This opened the way for labor effective remote MC generation.
- 2000 Remote MC production starts at Lancaster*
 - ◆ First version of 1st generation workflow manager, MCRunJob is used
- Prague, IN2P3 + other sites soon produce MC.
- 2002 TMBS pulled from FNAL for remote analysis to GridKa and IN2P3
 - ◆ FNAL central analysis capacity is undersized.
- “Regional analysis” discussions begin
 - ◆ Recognition in the collaboration that FNAL does not have to be only analysis site formalized in Offsite Analysis Task Force
 - ◆ SAR (Southern Analysis Region) is formed
- SAR develops mcfarm to as add-on to McRunJob, uses small sites effectively
 - ◆ In addition, supplied expertise to bring up farms in India and Brazil

Case Study: Global Production



- **2003-reprocessing with “p14”, planning at Beaune workshop in June 2003**
 - ◆ Takes about 2 months of prep for six weeks of processing
 - ◆ Carefully evaluated trade-offs
 - Could not use JIM
 - merged output at FNAL
 - ◆ Relied heavily on SAM bookkeeping.
 - ◆ 100 M of 500M reprocessed completed offsite.
 - ◆ NIKHEF used EGEE components, valuable learning experience.
 - ◆ Learned that certification time intensive as feeding the farms.
- **2004- preparation for “p17” reprocessing**
 - ◆ Goal is 1B events produced offsite in six months
 - ◆ Using JIM, merging at remote sites

Case Study: Global Processing

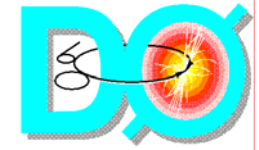
- 2004 offsite analysis tails off as FNAL resources improve.
- 2004 International Finance Committee recognizes computing as in-kind contribution based on level of service.
- 2005 Processing in progress
 - ◆ Lots of hard work by DO collaborators and site administrators
 - ◆ Working to improve efficiency
 - ◆ At this point, should finish reprocessing in about 6-8 months.
- In retrospect
 - ◆ People are excited by exciting projects.
 - Exciting computing projects can add to the physics.
 - Project leads decisions are respected—cohesive effort.
 - ◆ Value adding works and leveraged collaboration resources we couldn't have foreseen
 - We didn't value add for the analysis tools...
 - ◆ The devil is in the details—certification unexpectedly time consuming
 - ◆ Simultaneous development and operations is very stressful when effort is limited.

Summary



- *DO computing works quite well, and the offline software is meeting the basic needs.*
 - ◆ *Extra slides at the end.*
- *Good planning will lead to a good model. The DO model differs in detail, but basically followed the 1997 plan through changes in leadership and technology*
 - ◆ *Data format size underestimated, reco time SERIOUSLY underestimated, but lots of assumptions are correct (analysis needs effectively equal processing needs, etc).*
- *Planning is the best the hedge against the unexpected, particularly in managing budget shortfalls, can make informed trade-offs.*
- *Meeting the user's most basic needs in areas where they interact directly with computing system should be an extremely high priority.*
 - ◆ *Easier to find motivated people interested in computing to work on computing problems than it is to find motivated people interested in physics to work on computing problems.*

The Cuts



- Lots of little nips and tucks
- Some of the access pattern functionality of SAM was never implement because cheap disk made it less necessary
 - ◆ The Freight Train
 - ◆ Coherent PickEvents
- Various database projects were cut or descoped or are very, very late.
- Online streaming
- Run control parameter redux
 - ◆ Some use cases not addressed.

Hope springs eternal

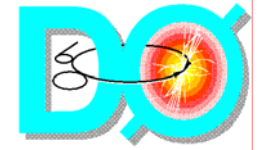


Algorithm software

- ◆ Reco is slow
- ◆ It takes a long time to cut a certified release
- ◆ Algorithm documentation
- ◆ Algorithm support and development
- ◆ MC tuning
- ◆ Putting PMCS into production

Attempting to address the shortcomings of Run I algorithm software motivated much of the work that went into software infrastructure (the framework, event data model, the run control parameter system, CVS and SoftRelTools). Had planned on a community of trained developers to help the less trained. Unfortunately, many people started from examples of bad code.

Implemented, but



- The features of the event data model, the framework
- Few users fully take advantage of SAM's bookkeeping facilities (but processing does)
- Online streaming (eventually?)
- Use of dCache write pools (eventually)
- A number of utilities were produced as part of the joint projects