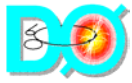# Computing Experience from CDF and D0

## Stephen Wolbers
## Fermilab
## CHEP2003, March 24, 2003

# Outline

- **Run 2 overview**

- **Computing experience/issues**

- **Mid-course corrections**

- **Future directions**

- **Conclusions**

**My viewpoint/bias: Deputy Head of FNAL CD 1997-2003, technical work on CDF production farms.**
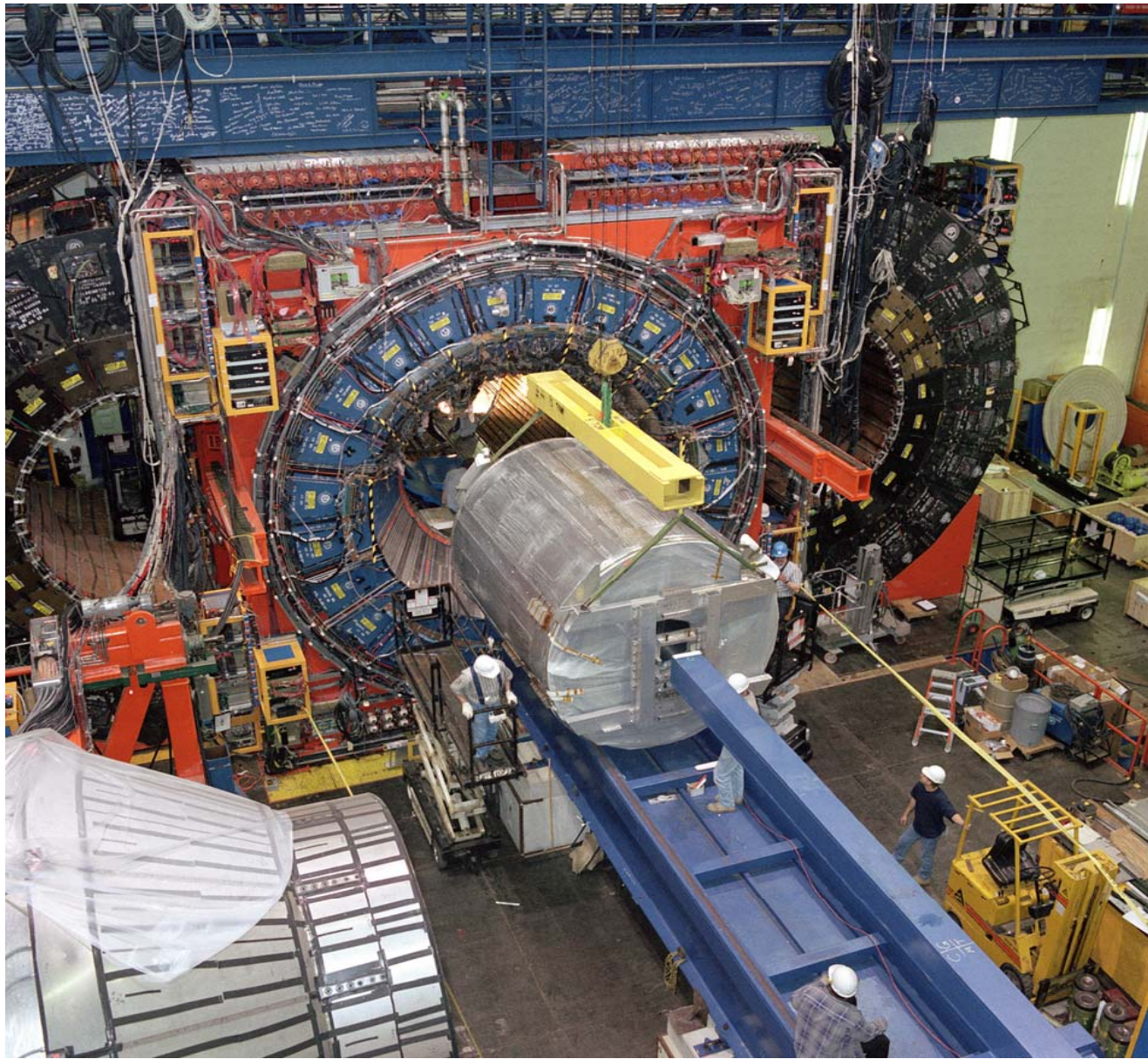
# Run 2 Overview

- Run 2 officially began March 1, 2001.
- Planning for Run 2 computing started many years ago and work has been continuous.
- Luminosity lower than expected, but it is increasing and will continue to increase.
- Both CDF and D0 have been taking data steadily at high rates limited primarily by DA capability. (Triggers are adjusted as luminosity increases.)
- Big detectors, large collaborations, many challenges.
- Computing is a big issue – essential for physics.
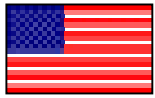- Run 2 has many years to go.

# The CDF Collaboration

## North America

3 Natl. Labs
28 Universities

2 Universities

12 countries

59 institutions

706 physicists

## Europe

1 Research Lab
6 Universities

1 University

4 Universities
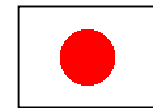
2 Research Labs

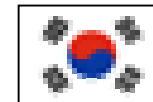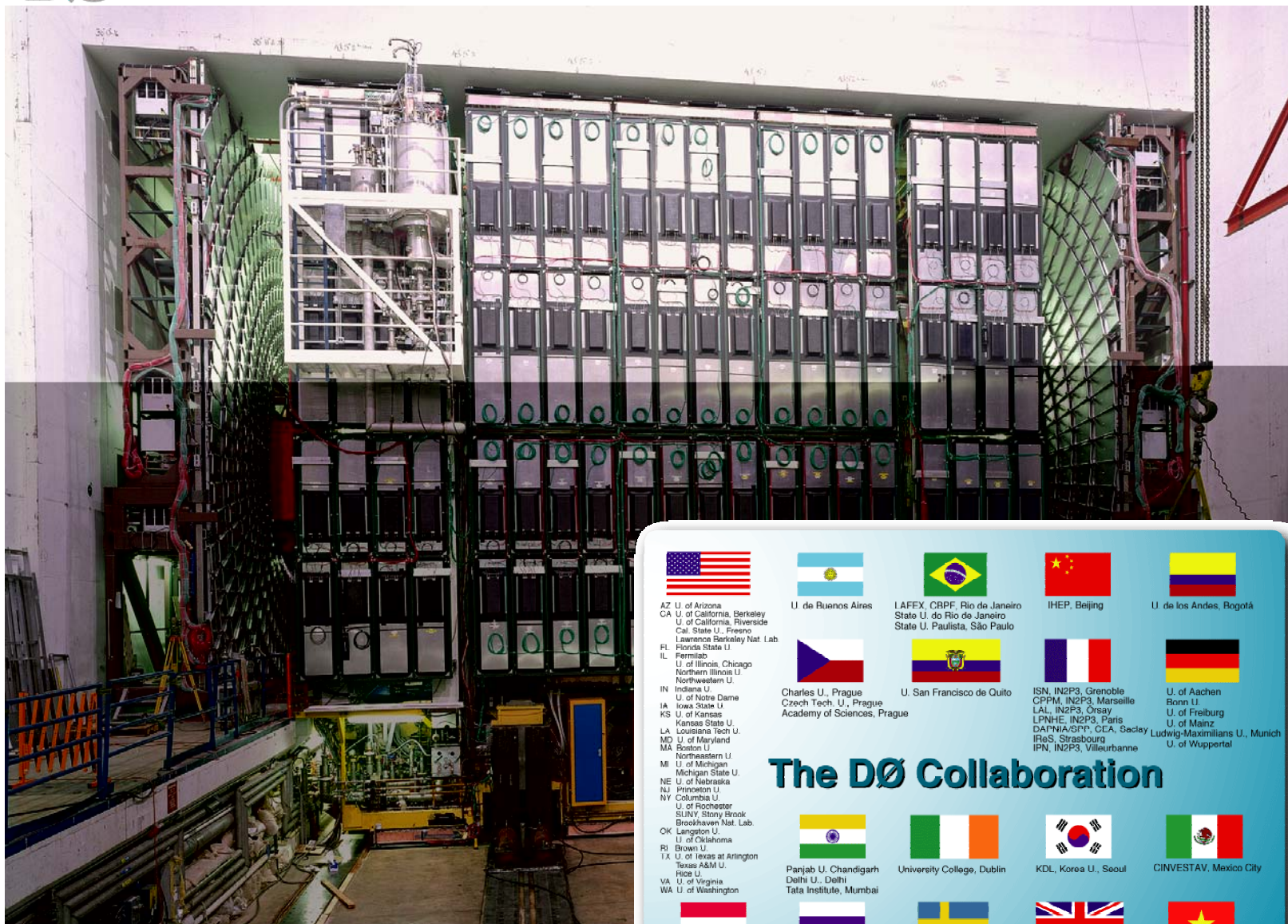1 University

1 University

1 University

## Asia

4 Universities
1 Research Lab

1 University

3 Universities

The DØ Collaboration

Stephen Wolbers

**Collider Run IIA Peak Luminosity**

Stephen Wolbers, CHEP2003

8

**Collider Run IIA Integrated Luminosity**

## CDF Raw Events



>180 pb$^{-1}$

500 Million events

# Run 2 Computing History

- **Early work/planning (goes back to 1995-1996):**
  - **Computing projects in both collaborations.**
  - **Computing Division efforts.**
  - Joint projects Computing Division/CDF/D0.
- **Planning and Reviews:**
  - DMNAG report (1997)
  - Software Needs Assessment (1997)
  - Von Ruden external reviews:
    - 1997(2), 1998, 1999(2)

# Short History (2)

- **The Run 2 joint computing project has so far:**
  - Designed hardware and software systems for Run 2 data storage, processing and analysis.
  - Defined and ran joint CD/D0/CDF projects.
  - Procured, installed, integrated and operated Run 2 offline computing systems.
  - Spent approximately $18M on equipment over 5 years (2-6-2-4-4).
- **Framework, code management, code distribution**
- **Reconstruction packages**
- **Online Systems**
- **Analysis Systems**

# What was successful?

- Joint Projects
- Data volume was more or less correct.
- Reconstruction farms.*
- Linux.
- Open source.
- ROOT.
- PCs.
- C++/C++ gurus.     Liz Sexton
- Networks & TCP/IP.*
- Mass Storage/Enstore.*    Don Petravick
- SAM data handling system.    Lee Lueking
- Offsite Monte Carlo production (D0).
- Commodity computing.
- Analysis Model.
- Code management/distribution.    Art Kreymer
- Reviews.*
- Moore's Law.*

* More Details will be given

# Farms

- **Farms (PC clusters) for event reconstruction and Monte Carlo are extremely successful.**

- **CPU power is plentiful and cheap.**

- **Networking is adequate, local resources are sufficient (memory and disk).**

- **The big issues are power and cooling and space!**

- **Maintaining these big systems is a non-trivial effort:**
  - **Hardware/OS/Networks/NIS/NFS.**
  - **Reconstruction Code.**
  - **Database and mass storage connections and tuning.**
  - **Operations – Software and coordination with collaboration.**
    - It is easy to mess this up!

# Farms – Run 1 vs Run 2

- 30 MHz UNIX WS
- 16-32 Mbyte memory
- Shared 10 Mbit
- 10 Mbyte executable
- 5-20 seconds/event
- 5-7.5 Hz
- 100 CPUs

- 1.7 GHz PC            (X60)
- 1-2 Gbyte memory      (X60)
- Switched 100 Mbit     (X100)
- >200 Mbyte executable (X20)
- 2-10 seconds/event    (X30)
- 50-75 Hz              (X10)
- 800 CPUs              (X8)

**A tremendous increase in capability**

# CDF and D0 Reconstruction

**Bottom Line: Keeping up with Raw data and Reprocessing**



D0



CDF



D0 CPU utilization

# Versions of code/CDF

| | | |
|---|---|---|
| I (Commissioning) | 3.11.0g | 9,775,297 |
| I (Commissioning) | 3.12.0 | 9,866,564 |
| II (1x8) | 3.14.0 | 2,539,217 |
| III (36x36) | 3.15.0c | 5,113,927 |
| III (36x36) | 3.16.0 | 5,113,927 |
| IVa (June, 2001) | 3.17.1 | 12,322,465 |
| IVb (June, 2001) | 3.18.0 | 21,993,586 |
| V (August, 2001) | 4.0.0i | 9,882,486 |
| V (August, 2001) | 4.1.0 | 38,857,497 |
| V (August, 2001) | 4.2.0a | 8,685,950 |
| VI (December, 2001) | 4.2.0b | 6,445,760 |
| VI (December, 2001) | 4.2.0c | 26,944,076 |
| VI (December, 2001) | 4.3.1 | 14,519,033 |
| VII (2002 data) | 4.3.1b | 10,369,500 |
| VII (2002 data) | 4.3.2 | 10,109,874 |
| VII (2002 data) | 4.3.2a | 38,690,220 |
| VII (2002 data) | 4.5.2 | 125,161,272 |
| VII (2002 data) | 4.8.0 | 16,572,810 |
| VII (2002 data)(through Jan,2003 shutdown) | 4.8.4a | 257,886,088 |
| TOTAL | | 630,849,549 |



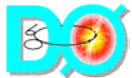**750 Million Events**

**Includes Reprocessing**

**CDF Farms**

**D0 Farms**

# Networks

- ## LAN
  - ### Very large and growing LAN for both experiments:
    - Connection of central systems.
    - Desktops.
    - Data movement drives the requirements ever higher.
  - ### A tremendous success.
  - ### Issues of scaling need to be solved.
    - Especially switch-to-switch.

# D0 Run-II Network Topology

# Wide Area Networking

- **This is an area which is becoming ever more important (widely distributed computing, grid, data exchange)**

- **ESNet upgrade to 622 Mbps December, 2002.**

- **Long-term increase in data rates have been seen.**

# Upgrades to WAN

- **Increased capacity for offsite connections is very important for Run 2 (and CMS and others).**
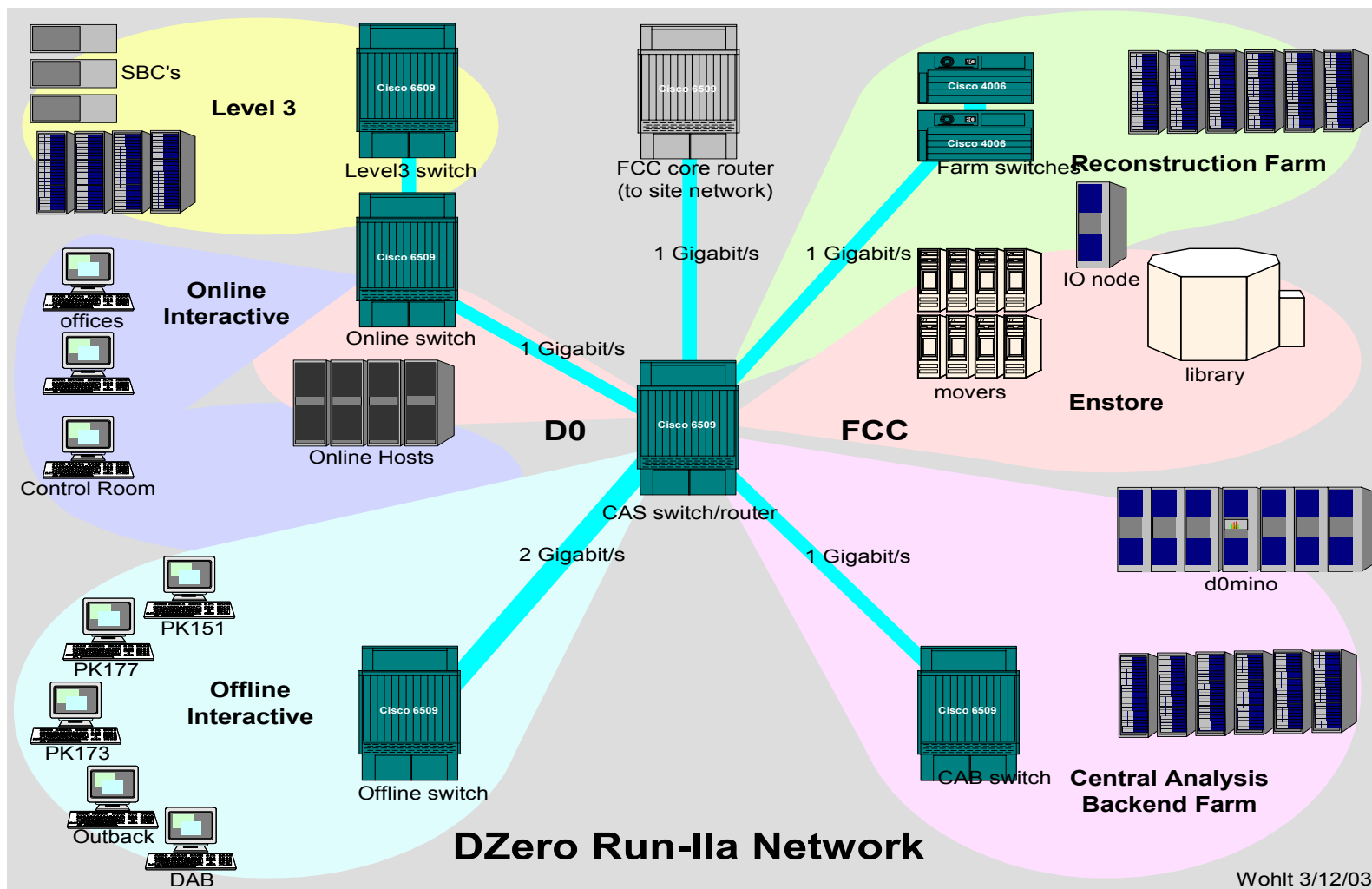
- **One possibility: Connect to Starlight (dark fiber leased from ComEd).**

- **Investigating many ways to increase bandwidth.**

Starlight: optical networking interconnection point downtown Chicago (*710 Lake Shore drive*)

Owned by Northwestern University

STARLIGHT

# Mass Storage

# Mass Storage

- **Mix of STK 9940a (60 Gbyte cartridges) and LTO 1 (100 Gbyte cartridges).**
  - **Allows competition – no single vendor.**
- **Currently migrating to STK 9940b (200 Gbyte cartridges) and investigating LTO 2.**
- **Software layer – Enstore.**    Don Petravick
- **Total Run 2 data on tape : 554 Tbyte.**
  - **Grows by ~2 TB/day.**
  - **Sometimes shrinks.**
    - Older farms output.



CDF Run II Data Logging March 2001 - March 2003

All (raw+processed) data
raw data only

Data volume stored on tapes [TB]

time [weeks]

some older produced data was
deleted to free up tapes

# Total Tape I/O per day



Stephen Wolbers, CHEP2003

Total Bytes Read Per Day (Plotted: Wed Jan 1 08:18:14 2003)

**20 TB/day**

**Disk I/O from CDF**
**Disk staging system**

dCache
Enstore

# Reviews

- **Reviews are good! (up to a limit)**
- **They allow for:**
  - **Organization and coordination of plans/thinking.**
  - **Collection of important planning information.**
  - **Outside/new look at computing issues.**
    - BaBaR, JLAB, CERN, DESY, other
  - **Access to new resources.**
  - **Attention of the Directorate/other higher level people.**

# Moore's Law

- **Moore's Law is essential to modern HEP computing capabilities.**
  - **We heavily rely on it.**
- **It is not a substitute for:**
  - **More efficient, faster code.**
    - 50% more CPU not the same as faster code, same amount of CPU
  - **Smaller datasizes.**
  - **Thinking before doing.**
- **The increased computing drives the science and vice-versa.**
- **Everyone would benefit from optimizing wherever possible.**
- **(Except those that give CHEP talks and need to show pictures of large amounts of computing stuff.)**

# What did not work (or needed change)?

- **Mass storage.***
- **SMP and Analysis computing.***
- **Scaling systems too quickly.**    *\* More Details*
- **Commodity PC and Disk.***
  - Hard to get and keep something that really works.
  - Linux kernels, large disk systems, RAID systems.
  - Disks.
    - Same issues as with PCs or most commodity equipment.
- **Windows NT/2000/Commercial Software.**
  - KAI (soon to be gone), analysis tools.
- **Fibrechannel/SAN.**

# What did not work (2)?

- **Data handling (CDF).***

*More Details*

- **Power and cooling of computer centers.**
- **Database performance.**
- **Reconstruction code speed.**
- **Lack of adequate monitoring.**
- **Procurement latency.**
- **wbs project management.**

# Mass Storage

- **8mm tape and flexible robots were chosen as Run 2 mass storage systems.**
  - Many reasons for this, all correct at the time.
- **Significant problems with performance, robustness and capabilities of 8mm tape and its integration into ADIC robots.**
- **Solution: STK/9940 and ADIC/LTO, coupled to the Enstore mass storage software (and Dcache).**
  - Rapid deployment, building modifications.
  - Successful major modification in mid-project.
- **Technology and competition gave us a solution.**

# Analysis Computing/SMP/PC

- **Analysis Computing in Run 2 is rapidly moving from big SMPs to PC based systems.**
  - Driven by cost, performance, capabilities.
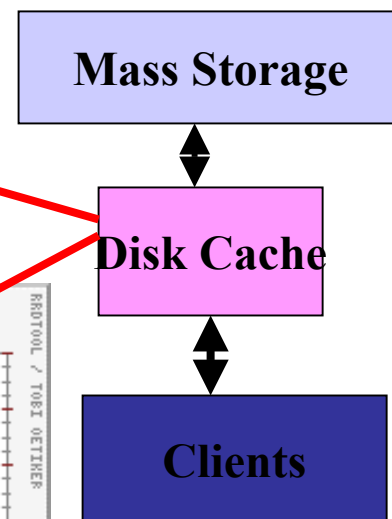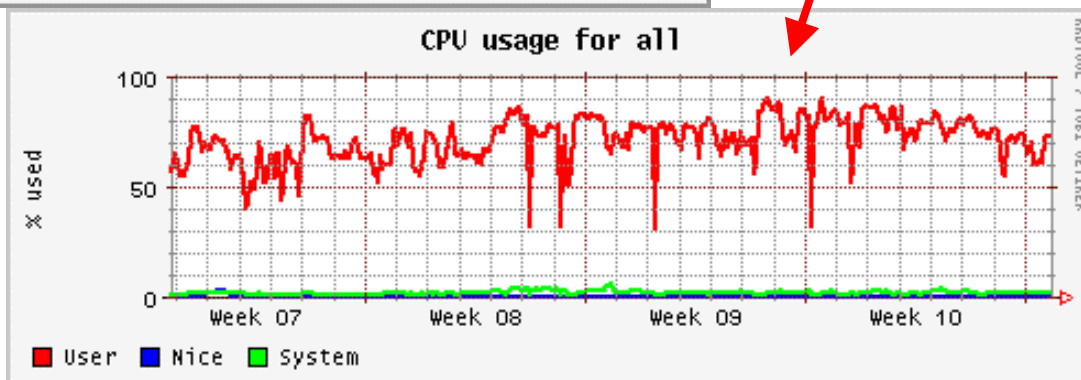  - SMPs remain as fileservers and common shared environments.

# Analysis Computing

- **Tough problem: Many clients, unpredictable behavior, short timescales (conferences), large peaks of load.**

- **Last problem to be addressed – DA, online, production, Monte Carlo all dealt with first.**

- **In addition, analysis computing normally competes to some extent with production activities:**
  - **Tapedrives**
  - **Network**
  - **Staging systems**

- **Analysis is distributed worldwide – should be integrated in the design.**

- **In the end the analysis computing is where the physics results happen – Very Important!**

# Moving a lot of data on CDF CAF



Network throughput for all

**600 Mbyte/s**

■ Received ■ Transmitted

Network throughput for fcdfdata077

■ Received ■ Transmitted

CPU usage for all

■ User ■ Nice ■ System

**Mass Storage**

**Disk Cache**

**Clients**

**Frank Wuerthwein**

Stephen Wolbers, CHEP2003

35

# Commodity Computing

- **Commodity computing is a challenge.**
  - **Much variety, constant change.**
    - Even within individual purchases.
- **Attempt to solve:**
  - **Common evaluation of equipment.**
  - **Common specification of configurations.**
- **Likely to remain an issue.**
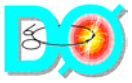- **This can be a huge effort!**

# Data handling/CDF

- **CDF data handling model became unsupportable.**
  - Direct attached tapedrives.
  - SMPs.
  - Fibrechannel/SAN.
  - Staging software.
- **Solution: move to common CD/CDF/D0 tools.**
  - Enstore/STK/9940.
  - Dcache.
  - PC analysis systems (CAF).
  - SAM/grid.
- **Effort to change was not small.**

  Gabriele Garzoglio, Fedor Ratnikov
  Rob Kennedy, Dmitry Litvintsev

  - Tape copying.
  - Software modifications.
  - System tuning.

# What was not expected?

- **Distributed computing, the grid.**
- **Really cheap PCs and disks and Linux.**
- **Long ramp up of the collider.**
- **Commissioning of the detectors and triggers.**
- **Duration of Run 2.**
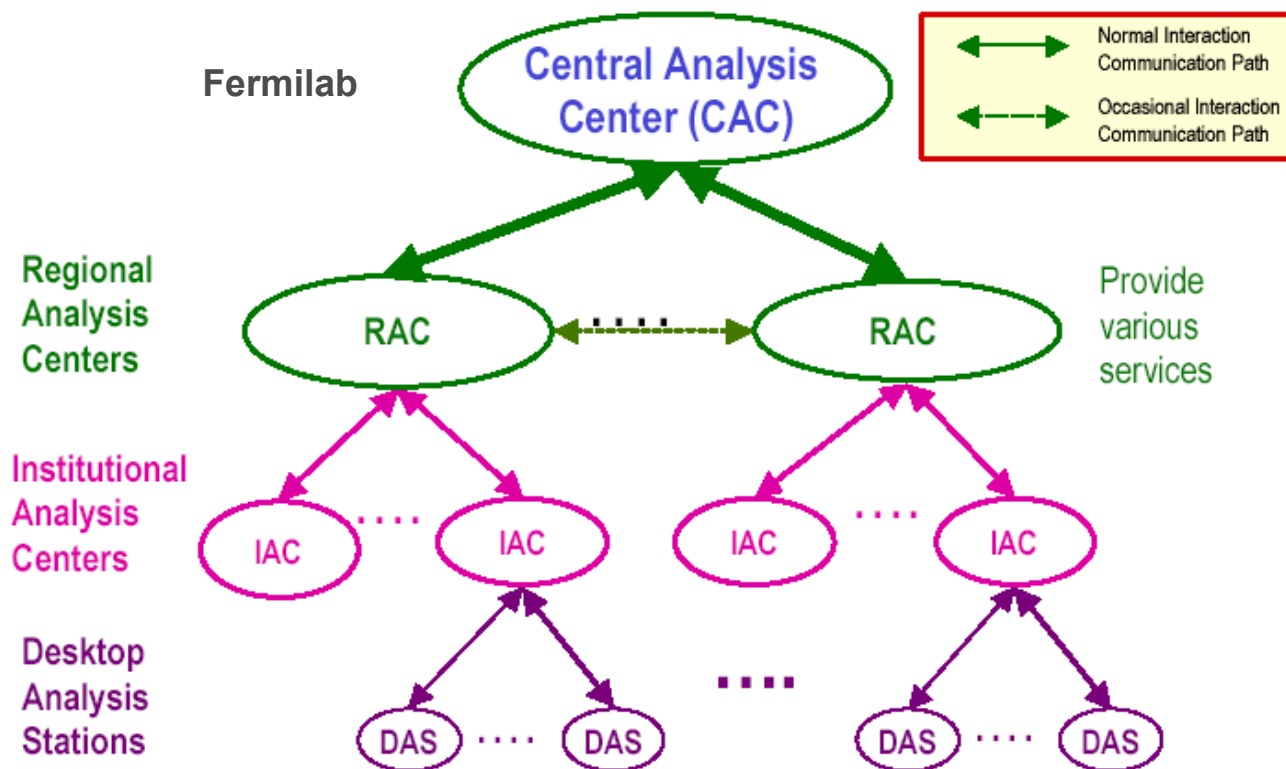- **Timescales for finishing software projects.**

# Distributed Computing and the grid

- CDF and D0 are both highly distributed collaborations, with many physicists and computing resources.

- Making use of that potential has been an issue from the start for D0 (Monte Carlo generation) and for CDF more recently.

- These efforts are going to grow.

- The SAM data handling system, used from the beginning by D0 and recently by CDF, is being used/modified for grid/distributed computing.

**Igor Terekov, Stefan Stonjek, Fedor Ratnikov, Lee Lueking**

# D0 regional analysis centers



## Proposed DØRAM Architecture

Fermilab

**Central Analysis Center (CAC)**

| | Normal Interaction Communication Path |
| | Occasional Interaction Communication Path |

Regional Analysis Centers — RAC ...... RAC — Provide various services

Institutional Analysis Centers — IAC .... IAC    IAC .... IAC

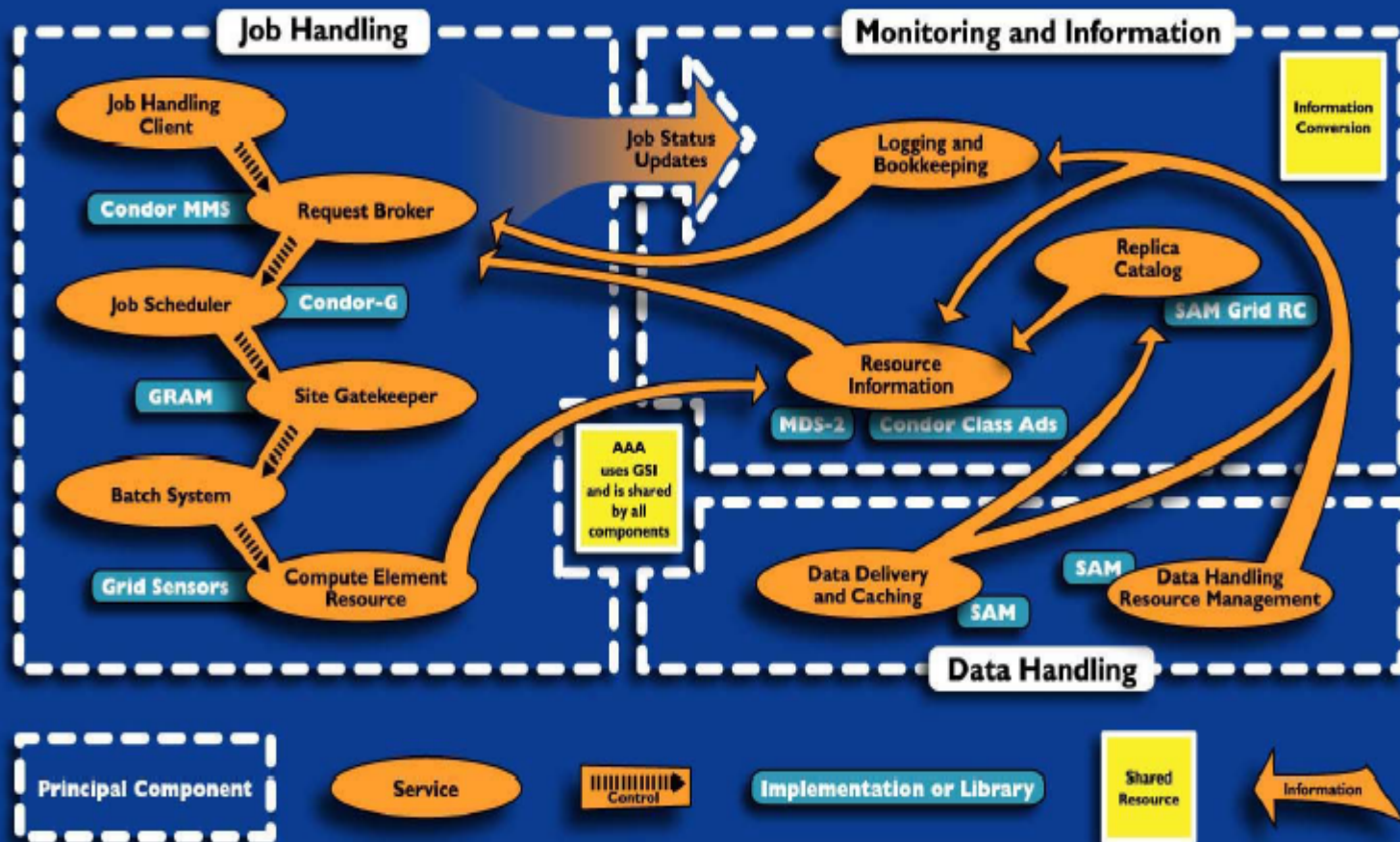Desktop Analysis Stations — DAS .... DAS  ....  DAS .... DAS
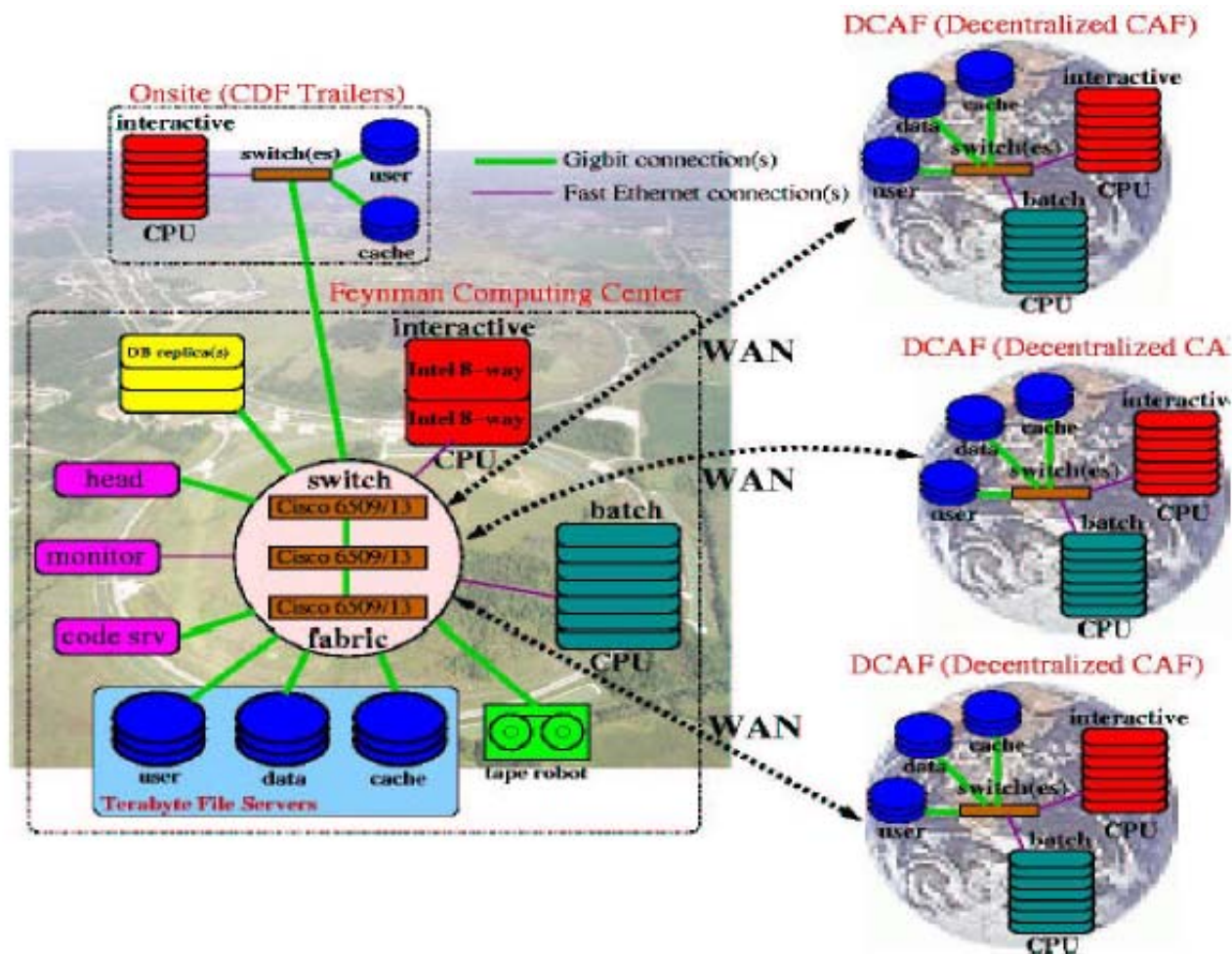
June 6, 2002

DØRAC Report
DØRACE Meeting, Jae Yu

3

SAM-Grid Architecture
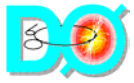
# Future CDF Direction

# Lessons

- **Plan for change, don't be afraid to make major modifications.**
- **Make realistic requirements estimates, taking into account experience as well as wishes.**
- **Optimize as much as possible.**
  - Hardware can make up for excesses.
  - But it is easier to run less stuff and/or get things done more quickly and easily and with less storage, memory, etc.
- **Don't buy too early.**
- **Test but be prepared for many surprises.**
- **Scale systems carefully and slowly.**
- **Be flexible, even late in the process.**
  - But in 10-20 year experiments it is not clear when "late" is.

# Lessons (2)

- Since all development (H/W and S/W) cannot be provided at one time, choose a path that gives necessary features and performance at any given time without impacting the overall system development (easy to say!)
  - This may require unoptimized systems at first, maturity coming later.
- Have a core group be responsible for infrastructure.
  - But be sure they are listening to the collaboration and others.
- Make good use of reviews.
  - Focus the project.
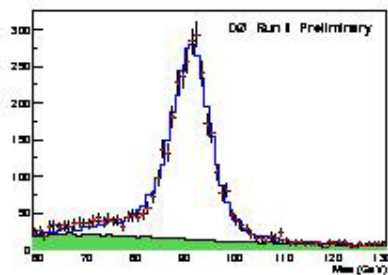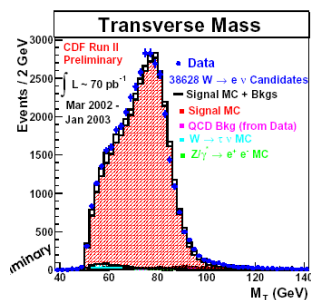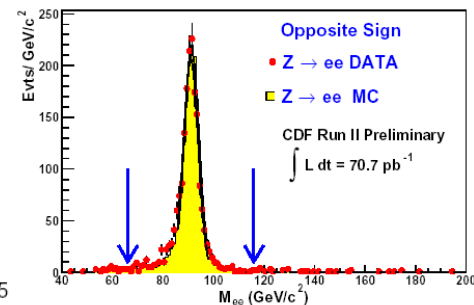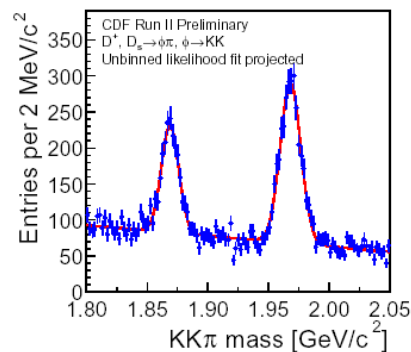  - Get attention of management.
  - Get resources.

# Lessons (3)

- **Joint projects are hard, but worth it.**
  - Require coordination with (many) other parties.
  - But the long-term support and features are worth it.
- **Two coordinators/leaders are better than one.**
  - Complementary strengths.
  - Better coverage.
  - Ideas can be bounced off of each other.
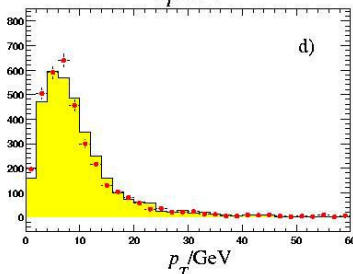- **Databases are very important.**
  - They connect everything.
  - Provide some reward for working on them.
  - Use freeware DB?

# Physics!

# On to the future

- **Run 2 is ramping up, has a long life.**
  - **Code maintenance.**
  - **Hardware upgrades and increments.**
  - **Technology changes.**
- **BaBar, RHIC, JLAB, LHC, others all facing similar issues.**
- **Grid is happening.**
- **The temptation (often the necessity) is to solve only our own problems.**
  - **It makes sense to try to align efforts to provide maximum effort and sharing.**
  - **It is somewhat difficult to align running experiments with LHC experiments.**

# Required Computing FY02 – FY05 (CDF)



Fiscal Year: 05, 04, 03, 02

Lum (fb⁻¹) 4.1: 1.6, 1.3, 0.9, 0.3

Batch CPU (THz) 4.7: 1.8, 1.4, 1.0, 0.5

Farm CPU (THz) 1.3: 0.54, 0.33, 0.37

Static Disk (TB) 540: 200, 160, 98, 82

Read Cache (TB) 160: 60, 46, 28, 26

Write Cache (TB) 46

Disk I/O (GB/s) 4.9: 1.8, 1.4, 0.9, 0.8

Archive I/O (GB/s) 0.48

Archive Volume (PB) 1.7: 0.6, 0.4, 0.4, 0.3

# Summary

- **Run 2 Computing is a great success.**
- **This was a large effort, involving many people over many years.**
- **The data has been processed and analyzed quickly to produce physics results.**
- **But the computing systems aren't perfect:**
  - **Entire systems were replaced.**
  - **Development had to (and has to) be consistent with data-taking and data-analysis needs.**
- **There is a huge amount of data on the way.**
  - **The luminosity continues to increase.**
  - **The data rate will increase in Run 2b.**
  - **The physics demands it and the technology allows it.**

# Thanks

I wish to thank many many people who worked on Run 2 computing over the years. It was truly a large and distributed effort.

I also wish to thank those who gave me suggestions and ideas for the talk – Wyatt Merritt, Amber Boehnlein, Heidi Schellman, Liz Buckley-Geer, Rob Harris, Frank Wuerthwein, Liz Sexton, Don Petravick, Ruth Pordes, Matthias Kasemann + many others who have worked on Run 2 over the past many years.