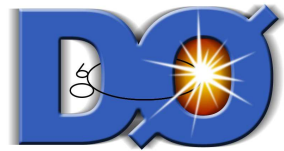


DØ Distributed Computing



Daniel Wicke
(DØ, Bergische Universität Wuppertal)



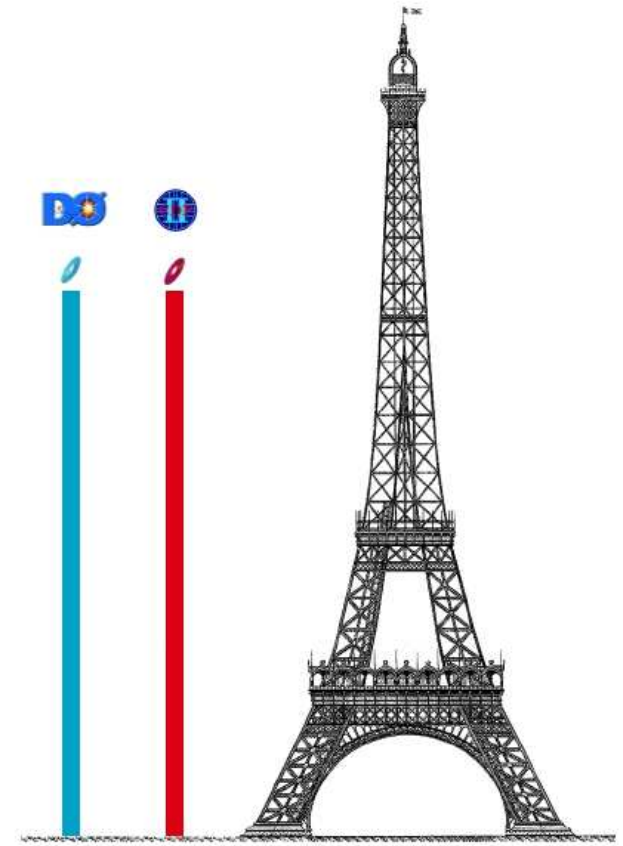
Outline

- Introduction: Computing Demands
- Computing Concepts
(Workflow Management, Local Resource Optimisation,
Distribution of Data and Jobs, GRID)
- Summary

Introduction

Computing Demands

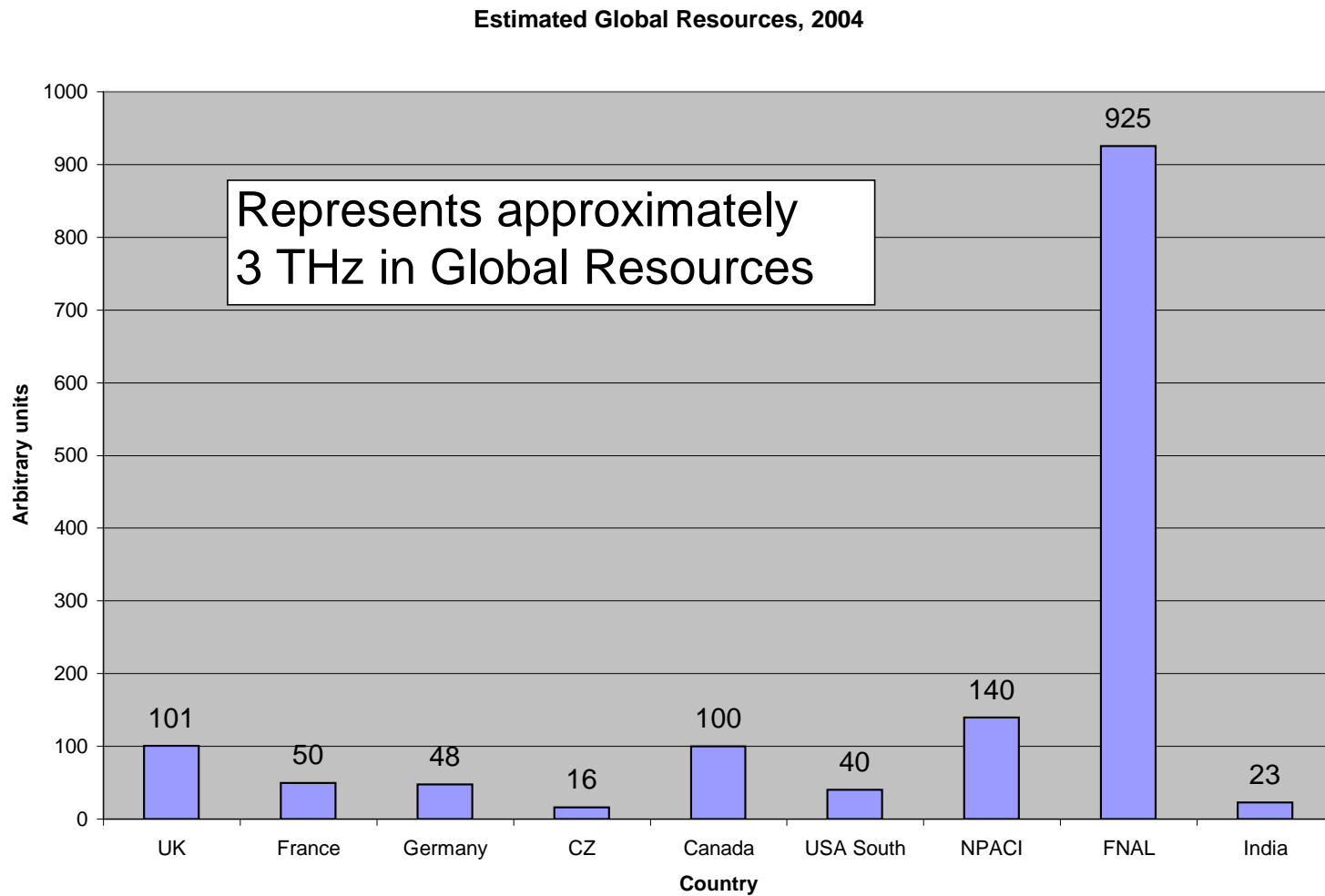
- Top and Higgs are difficult to recognise.
- To find sufficiently many top- and Higgs-events a large number of reactions must be recorded.
- CDF and DØ record 500GB/day each.
- Another 1100GB/day come from processing the raw data.
- This sums up to 200m of CDs per year.



Providing facilities to analyse these data is a big challenge.

Estimated CPU need in 2004: 1.4MSpecInt2000 (3000 1GHz PIII Computers).

DØ computing resources

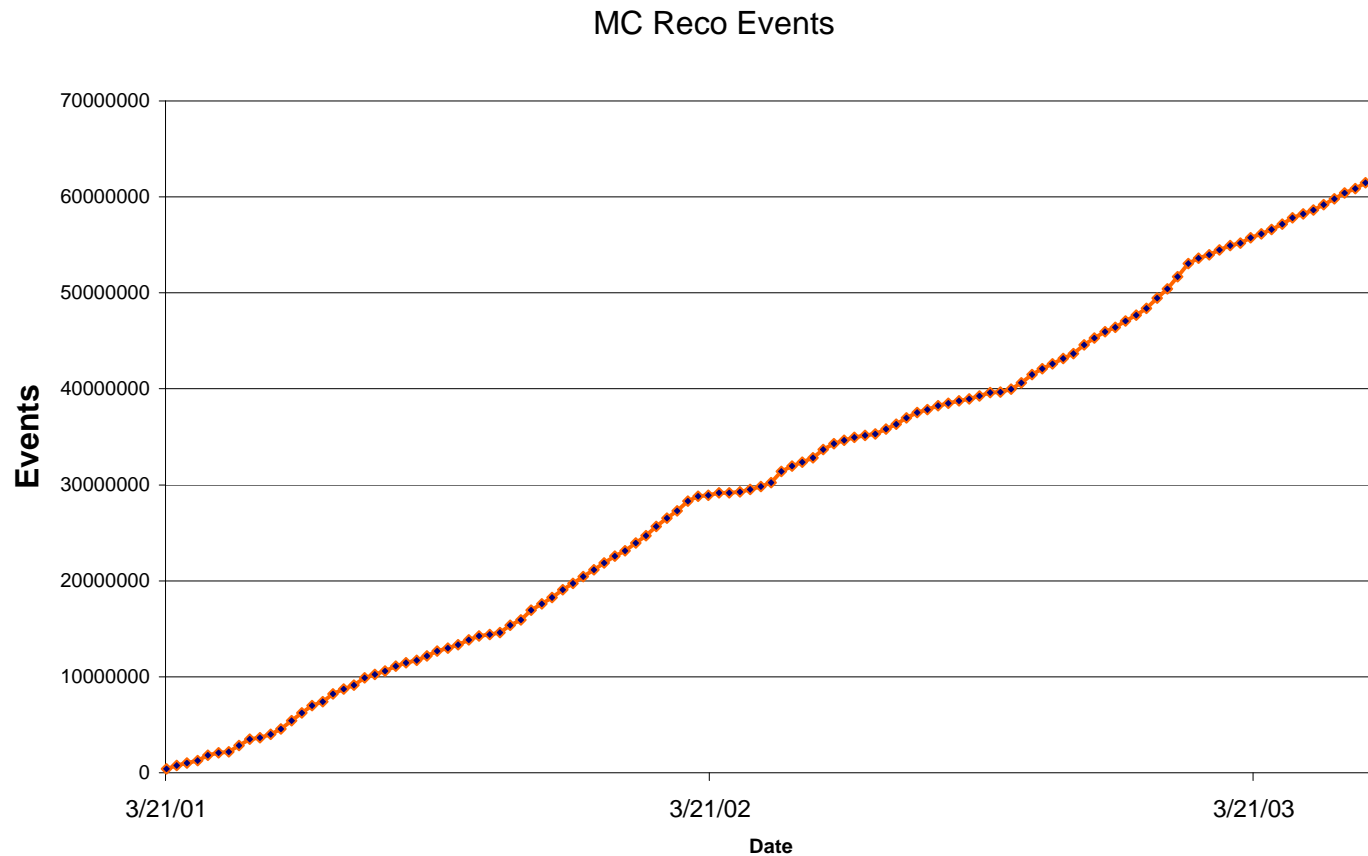


About half of the resources are located at remote sites.

Major computing tasks in DØ

- Monte Carlo Production.
 - Self-contained.
 - CPU intense with relatively little IO.
 - Simplest.
- User analysis jobs.
 - Requires access to DØ-data.
 - Mostly IO intense.
 - Chaotic by nature.
- Secondary reprocessing of data.
 - Requires access to DØ-data.
 - High CPU and IO demands.
 - Bookkeeping for special binaries required.

MC Production



- All simulated events are produced at remote sites
- The production has been stable over the last 2 years.

Management of Large Batches of Jobs

Problem

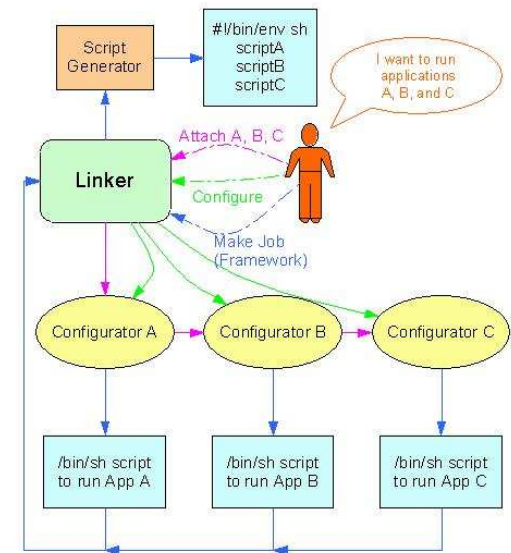
Handling MC production chains with ever changing program versions/parameters is very **person power intensive**.

Workflow management: Runjob

- handles chains of executables
- passes output of one program as input to a following
- track metadata
- allows for automatic parallelisation of jobs
- separates calling details from organisation

Base of improvements

Automate linking for several programs into a single job based on a job description (configurator).



Future Runjob developments in collaboration with CMS.

Data Analysis

Optimised Use of Local Resources

Problem 1

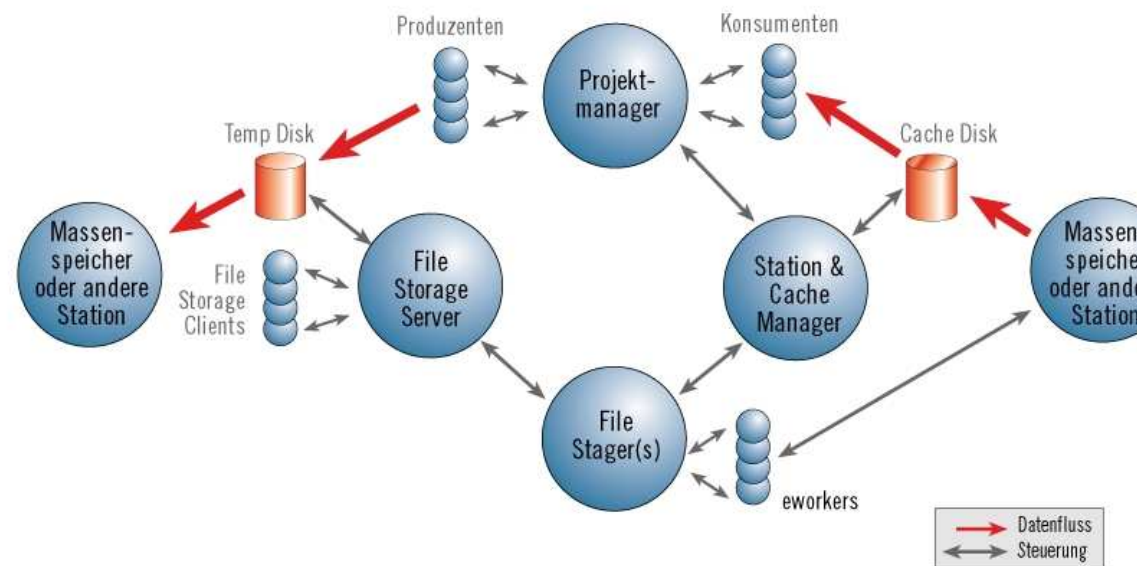
It is impossible to store all data on disks:
Tape access is the major bottleneck.

Optimisation with SAM

- Don't loop through file lists.
- Request **datasets**.
- The order in which files corresponding to a dataset are processed may change.
- The system optimises the order to minimise tape access and tape mounts.

Base of improvements

The order of events in the dataset has no meaning.



Worldwide Distribution of Data

Problem 2

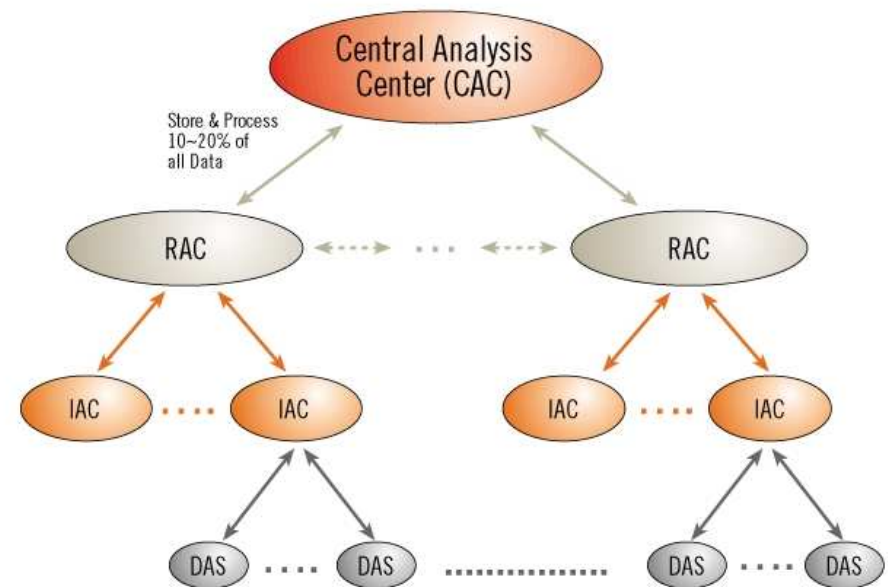
For the most frequently used data the I/O rate to disks limits performance.

Regional Analysis Centers (RACs)

- Allow full physics analysis.
 - ⇒ Hold all Thumbnails (microDSTs).
 - ⇒ Provide computing power to process these.
- Allow distributed reprocessing.
 - ⇒ Hold 10% of full DST.
- Serve institutes.

Base of improvements

Setup copies of these most frequently used data.



As such a major building block for the DØ Grid.

Establishing the RAC Concept: Setup of Prototype

Goals

- Proof of principle for RAC concept.
- Check needed resources and their scalability (esp. person power and network).
- Check DØ-software compatibility.

Specifications

The prototype should implement the following major features of a full RAC:

- Continuous and immediate transport of thumbnails from Fermilab to the prototypes disks.
- Fetching files available at the prototype from associated institutes.
- Automatic installation of DØ-software updates.

Prototype Location



GridKa: Grid Computing Centre Karlsruhe

- established in 2002 at FZK.
- German centre for Grid development.
- regional data and computing centre.
- 8 HEP experiments:
Alice, Atlas, Babar, CDF, CMS,
Compass, DØ and LHCb.
- 5 (now 6) DØ institutes in Germany associated

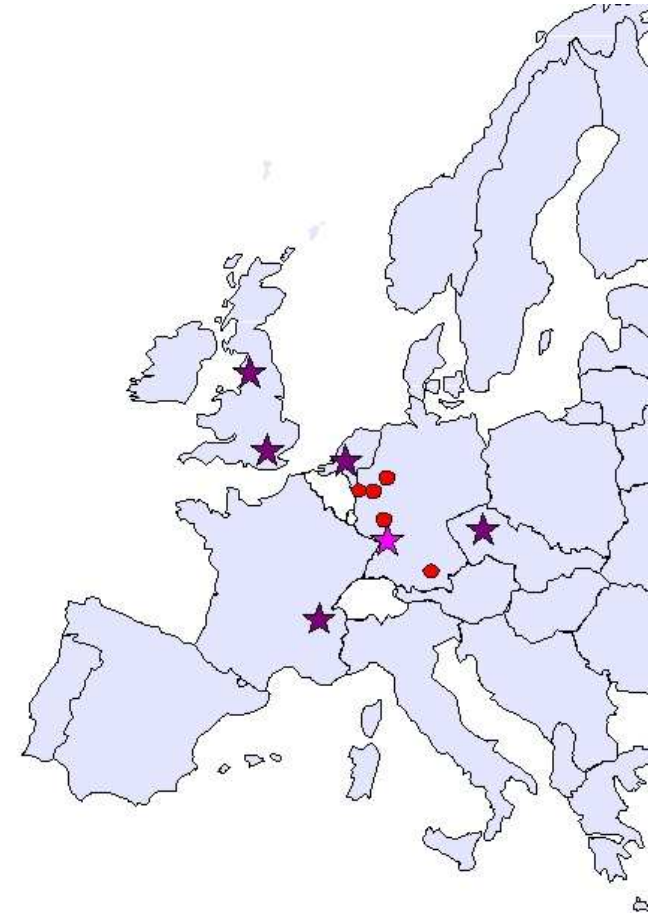
Resources

160 dual CPU nodes, 44 TB disks, >100 TB tape

Other (european) potential RACs

London, Lyon, Manchester, Nikhef, Prague

Some established for MC production since long.

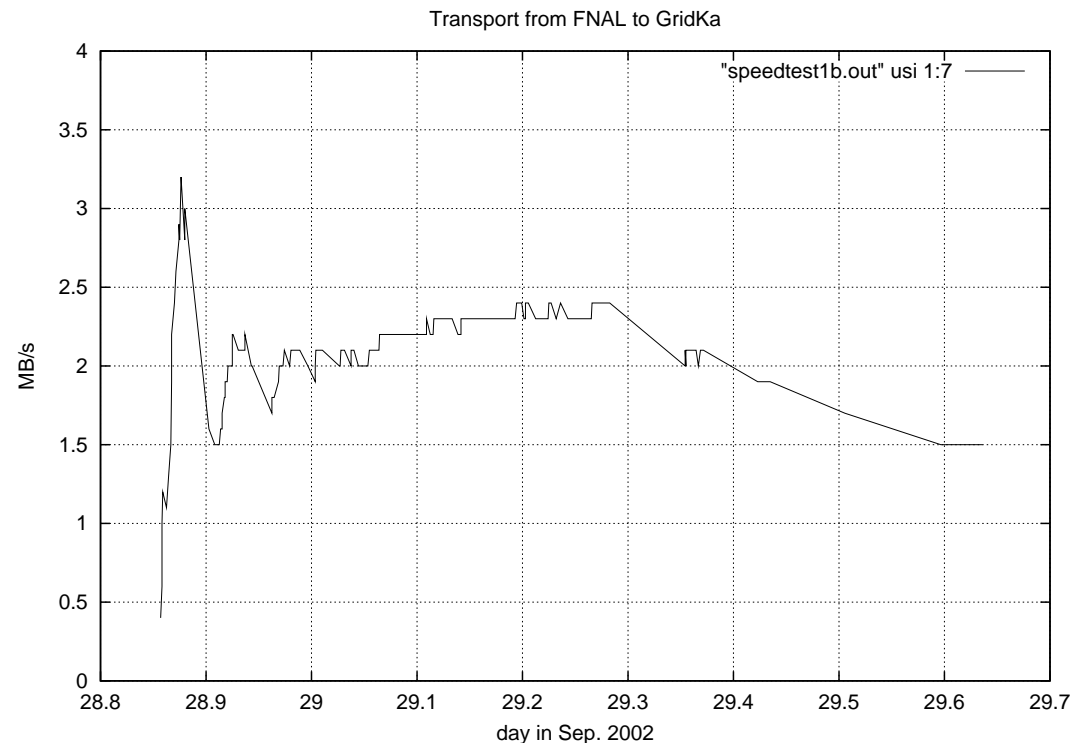


Results: Transfer Speed FNAL – GridKa

Integrated size of arriving files over time:

2–3MB/s = 15–25MBit/s
(averaged over 7 hours).

larger gaps in the transport decreased effective speed thereafter.



⇒ (At that time) limited by FZK connection (32MBit/s)
Currently upto 3.3MB/s (or 8TB/month) is observed.

Internet bandwidth to FNAL is sufficient for our current demand.

Software Problems Lead to Person Power Problems

DØ software originally written for the FNAL computing environment and still tested in this environment

Deploying this software at remote sites

- requires site specific adaptations
(hardware setup can't be changed everywhere, needs repetition for each new version)
- requires site specific configurations
(sometimes quite fancy and possibly fragile)
- introduces correlations between **all** sites.

Person power becomes a problem when running many systems

Example: Unfriendly correlations

The launch of a new linux cluster (CAB) at FNAL induced unexpected problems at GridKa:

Symptom

- Large number of transfer errors/undelivered files.

Causes

- GridKa SAM-station tried to get files directly from CAB.
- This failed because the CAB-nodes aren't known to our firewall.

Solutions

- Opening the firewall isn't feasible.
- SAM configured to route files from CAB through central-router.

Person power required to run the systems needs to be reduced!

⇒ Aim to remove/avoid dependencies between sites:
Think globally, avoid site specific paths, code and libraries.

Acceptance by the Users

Physics results

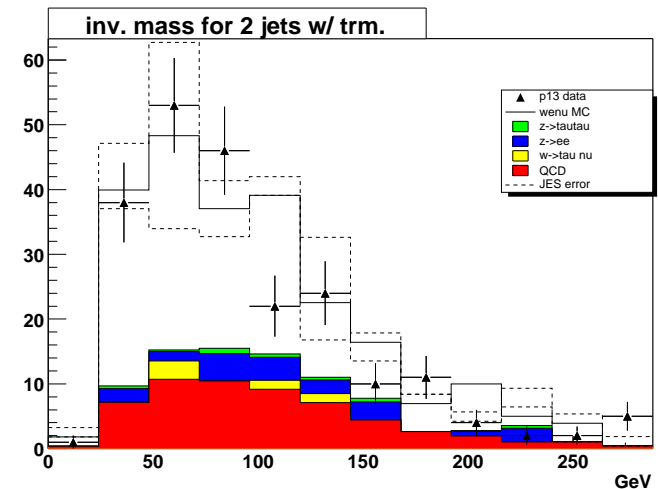
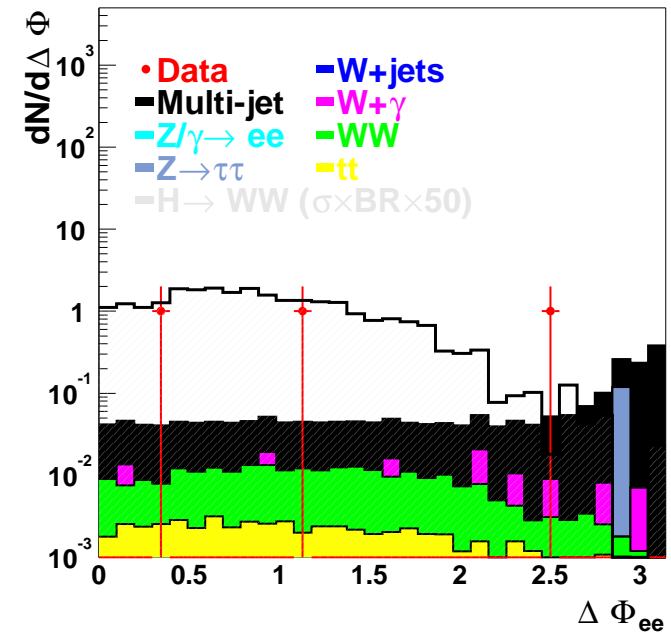
- $H \longrightarrow W^+W^-$
- W +jets

Analyses produced results for Moriond 2003.

MC Productions

- 1.2 million MC events for the top group
- Vecbos vs. Alpgen comparisons
- ...

January 2003: Concept was accepted by the users and has proven its value.



Secondary Reprocessing of Data

To achieve full physics return for Moriond data processed by older versions shall be rereconstructed with the greatly improved new event reconstruction software.

Organisation

- Central production farm can't handle the extra load.
- Large portions of the will be done remotely.
- Still not enough: Preselection or prioritisation of events/runs is needed.

Status

- Several scripts needed adaptation and are being tested.
- Reco-version without DB access needs to be certified.
- Physics groups are working on the preselection/prioritisation criteria.

Time-line

- Start of processing on the farm: Next week.
- Start of remote processing: October.

A GRID for DØ

Problem 3

Network, CPU and disk space availability will be constantly changing. Communication lines need to be configured explicitly.

Base of improvements

Dynamic adaptation to actual situation.

Optimised use of globally distributed resources: The GRID

- Retrieve data from remote disks with best availability.
- Submit jobs to centres for which they're best suited.
- Base these decisions on **current** status of the systems.

Requires

- Common protocols for data and status information exchange.
- The GRID

Summary

- DØ computing strategy is driven by physics needs.
- Computing demands are met by using
 - automated workflow and metadata management (runjob)
 - distributed data access with user-level bookkeeping and local optimisation (SAM)
 - global distribution of data (Regional Analysis Centres)
- Moving towards GRID technologies we hope to
 - reduce the person power needed to run the clusters
 - add global resource optimisation

Outlook

- Perform distributed data rereconstruction.
- GRID requirements will be tested in a running experiment.
- Evolve towards full GRID standards.