# CDF long term data preservation
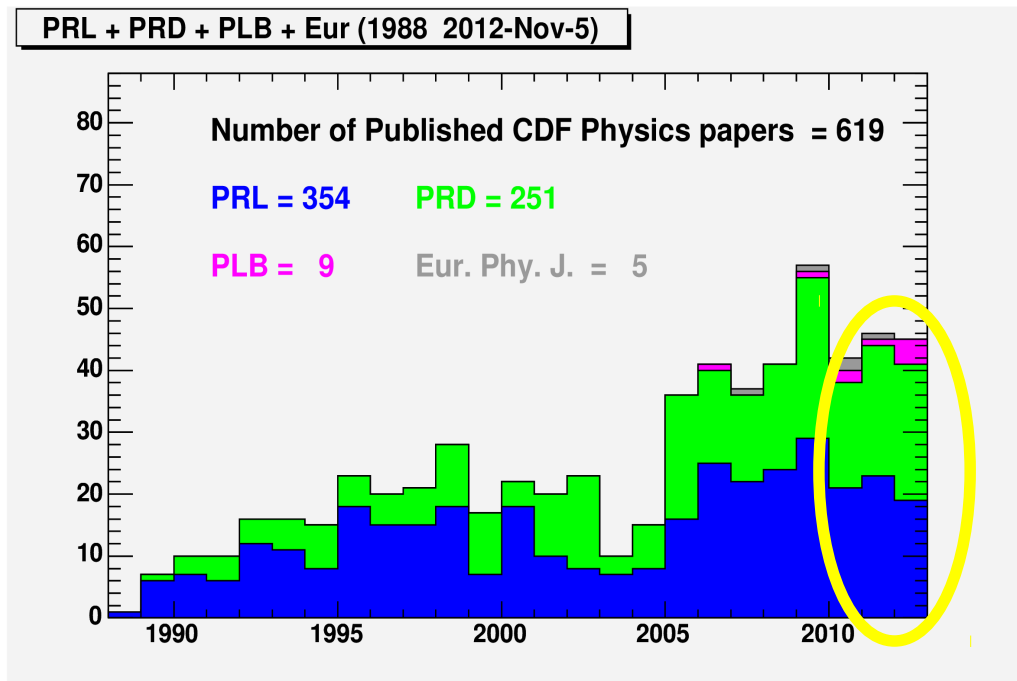## - November 19, 2012 -

**S. Amerio**

(INFN Padova)

on behalf of CDF data preservation task force

More than one year since the end of the operations.

**PRL + PRD + PLB + Eur (1988  2012-Nov-5)**

Number of Published CDF Physics papers  = 619

PRL = 354        PRD = 251

PLB =  9        Eur. Phy. J.  =  5

2011: 46 papers
2012: 40 papers

CDF analysis plan:
- complete and publish current ongoing analysis with the full data sample
- Analysis with legacy potential:
  • QCD and EWK measurements @ 2 TeV
  • QCD analysis of the energy scan
  • t-ttbar and b-bbar asymmetry
  • W boson and top quark mass
  • ….

*A lot of activity foreseen in the next few years.*
*But Tevatron data will remain unique and of great interest in the long term future.*

# CDF Long Term Data Preservation Project

*Goal: preserve CDF data and analysis capability in the long term future (> 10 years from now)*

We aim at Level 4 preservation: *full analysis chain* capability
- Data and technology to access it
- Analysis code
- Computing resources
- Knowledge!

A CDF long term data preservation task force is official since June 2012.

From our charge:
*"... Production of physics results should not be more difficult in the future than it is at present despite the expected diminishing support and reduced availability of expert advice...."*

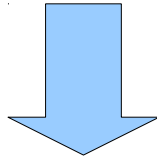Our target: physicists from CDF collaboration (and beyond !)

Current members:

- Bo Jayatilaka *(coordinator), S.Amerio (co-coordinator)*
- *Ray Culbertson*
- *Rick Snider*
- *Steve Wolbers*
- *Tyler Parsons*
- *John Strologas*
- *Donatella Torretta*
- *Physics groups conveners are ex-officio members*
- *Physics groups representatives:*
    - **B:** Satyajit Behari
    - **QCD+EWK***: Niccolo Moggi
    - **Higgs**: Craig Group
    - **Top**:  Yen-Chu Chen (all leptonic analysis), Jon Wilson (semi-leptonic analysis)

An intense summer...

- Regular meetings every two weeks.
- Mid-August:  Joint CDF/D0 meeting
- End of August: Report to Fermilab management

Report of CDF Data Preservation task force

CDF note CDF/DOC/CDF/PUBLIC/10922

First report released on Sept.17th.

In the report we identified requirements
and possibile solutions for:
- *data access*
- *software preservation*
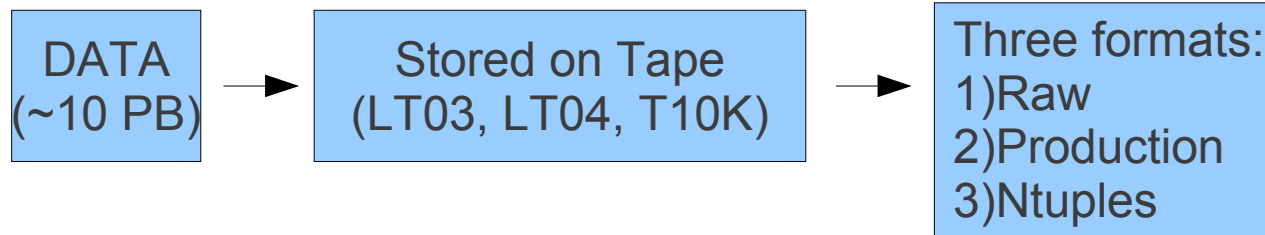- *job submission*
- *documentation*

S. Amerio[1], S. Behari[3], J. Boyd[3], Y.C. Chen[2,3], R. Culbertson[3], C. Group[4], S.Lammel[3], N. Moggi[5], R. Snider[3], J. Strologas[3], T. Parsons[3], D. Torretta[3], S. Wolbers[3]

[1] INFN Padova (Italy), [2] Institute of Physics, Academia Sinica (Taiwan), [3] Fermilab (USA), [4] University of Virginia (USA), [5] University of Bologna (Italy)
(Dated: September 17, 2012)

We are designing the project looking for
• The best compromise between available manpower from the collaboration
and  essential needs for analysis in the long term future
• Strict collaboration with D0 and FNAL Computing Sector.
• Feedback from DPHEP and other experiments

```
┌──────────────┐        ┌──────────────────────┐        ┌─────────────────────┐
│ DATA         │   →    │  Stored on Tape      │   →    │ Three formats:      │
│ (~10 PB)     │        │  (LT03, LT04, T10K)  │        │ 1)Raw               │
│              │        │                      │        │ 2)Production        │
└──────────────┘        └──────────────────────┘        │ 3)Ntuples           │
                                                         └─────────────────────┘
```

Physical tape
• Data will require to be regularly migrated to new tape technology
• Complete migration of LT03/LT04 data to T10K in the next two years.

Data redundancy
• Production data contains a copy of raw data
• A project to have a complete copy of CDF data offsite is being developed in collaboration with INFN

CDF Italian collaboration, together with CNAF computing center and FNAL, is developing a project to preserve at CNAF a copy of CDF data (raw and ntuples) and the analysis capabilities in the long term future.

**GARR network**

- Use GARR network (Italian R&E network) for the copy
- Adapt the current copy mechanism between FNAL and CNAF
- First feasibility tests performed during the summer → successful copy retrieval of a small dataset to CNAF tape system and retrieval.

- The project got a first approval by INFN at the end of September; final approval is expected in December.

*Data handling system* based on

- **SAM (Sequential Access via Metadata)**
  - developed at Fermilab, used by D0 and CDF, it closely integrates with Fermilab Enstore tape manager and other data handling systems (e.g. dCache)
- **dCache**: it fetches files requested by the users and stores them on a distributed pool of disk servers (800 TB) for the user to access over the network.

Data access in the long term future:

We have to maintain SAM; two options under study:

1. Keep the current system as it is now

2. Update to the new SAM being developed for the Fermilab experiments

Option 2. requires more effort in the short time, but will ensure support in the long term future.

*Databases:*

**Oracle** for

- run condition, configuration, trigger, luminosity, alignement, calibration →

*necessary for new MC generation;*

- data metadata (dataset, fileset, file and run section information) → *essential part*

*of SAM access method*

Long term future:

We prefer to keep Oracle; not enough manpower to migrate to a different DB and

perform all the validation.

FNAL will keep using Oracle in the future; we may need to migrate to new Oracle

versions in the future.

- CDF-supported software is archived as a set of packages in a **CVS repository:**

  - Fermilab is committed to preserve the content of thre repository indefinitely

  - Migration to SVN under investigation

  - ALL analysis code has to be archived (some still on user's desktops)

- Build configuration toolkit based upon **SoftRelTools** and **ups** products

supported by Fermilab.

- **SL5; new SL6 version in 2013.**

  - Likely the port to SL6 will be the last

  - We are considering **virtualization** to preserve the existing build and run-time

    platforms

- A well documented and tested **validation procedure** is needed

• CDF *Central Analysis Farm code (CAF)* provides the users with a uniform interface to resources on different Grid sites

**Job submission headnode**

CDF Central Analysis Farm code:
•Submitter
•Monitor
•Mailer

User pool (Condor)

**Factory**

Grid site

CDF Frontend

Grid site

• Three *portals* to access computing resources:
• CDFGrid → FNAL
• NamGrid → OSG
• Eurogrid → Tier1 @ CNAF and LCG

Based on *glideinWMS* workload management system (batch system = *Condor*)

In the near term future (< 5 years) the current system can be maintained with minimum effort.

Realistically, in the long term future access to CDF data will be greatly reduced; we need to move to a job submission system simple, flexible and easy to support.

Two options under study:

1) Development of a system which aligns with and integrates with the job submission system  being developed for the new Fermilab experiments:
- We will take full advantage of developments and support.
- This solution requires consistent effort to adapt the current system.
- Security issues need to be addressed.

2)  Build a virtualized system which can run in isolation
- No security issues.
- It requires dedicated support and maintenance.
- It may not be justified by the CDF analysis load.

**WNoDeS** → **Worker Nodes on Demand Service**

Virtualization architecture which provides transparent user interfaces to Grid, Cloud and local access to resources

In **production at several Italian centers, including the INFN Tier-1 since November 2009** (Currently managing about 2000 on demand Virtual machines there)

**New feature under development: Dynamic virtual networks:**
dynamic instantiation of **private** VLANs and address assignement for VM isolation.

*In the long term future at CNAF: CDF services and analysis computing resources can be instantiated on demand on pre-packaged VMs in a controlled environment.*

We need to carefully **preserve and re-organize** as much documentation as possible

- *Internal webpages*
- *Webtalks pages*
- *CDF notes archive*
- *Logbooks*
- *Twiki pages*

CDF-1

**Parameters of Colliding Beam Detector**

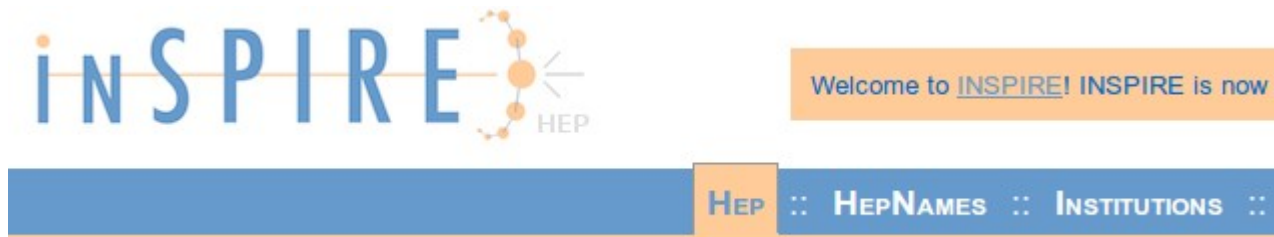R. Diebold, A. Tollestrup, T. Collins, S. Ecklund, J.K. Walker

23 Jan '78

Constraints

1. Kissing scheme for beams.

2. Maximum energy of MR/ED beams = 150/1000 GeV for pp and 1000/1000 for $\bar{p}$

3. Conventional magnets for the MR normal operation.

4. Low β achieved without loss of part of the 50m long straight.

   (Note: R. Diebold paper #1 of Summer Study assumes 46m is free.)

1000 notes from the Berkeley collection *scanned and uploaded* to our archive.

Thanks to Lina Galtieri, Stephanie Schuler and Barb Hehner!

- 11000 internal notes, from 1978!
- We will archive them in Inspire
- Still in the testing phase: a first set of notes will be uploaded in the next month
- If no major issues, all notes arvchived by next Spring.

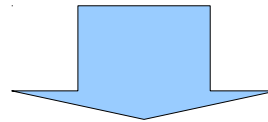Many thanks to Zaven Akopov, Heath O'Connell and  Alan Jonckheere.

CDF website **www-cdf.fnal.gov**
• Public section
• Private section
  • Online
  • Computing
  • Physics (with sub-pages for each physics groups)
  • Organization (CDF organization, meetings, webtalks, internal notes, ...)

CDF online webserver **www-cdfonline.fnal.gov**
• Data acquisition and sub-detector e-logs, details about detector systems, daq and trigger operations

Re-organize all the essential information in a new web-site

*Target*: a physicist, non necessarily from CDF
*Goal*: he/she has to be able to perform an analysis from the very beginning
*Timeline*: 2 years
*How*: 1) identifiy representatives for the different sections 2) identify well defined tasks that should be considered as service tasks

The CDF data preservation task force is official since June 2012.

*First report released In September 2012.*

We  have identified necessary requirements to preserve in the long term future CDF analysis capabilities.

*Preservation of the documentation has already started (high priority, given the natural reduction of the collaboration).*

CDF/D0 requirements are being discussed with Fermilab Computing Sector experts; both task forces will be part of the Fermilab data preservation project team (details in Steve's talk).

- Backup -