

CDF Computing and Analysis Model

January 26, 2009

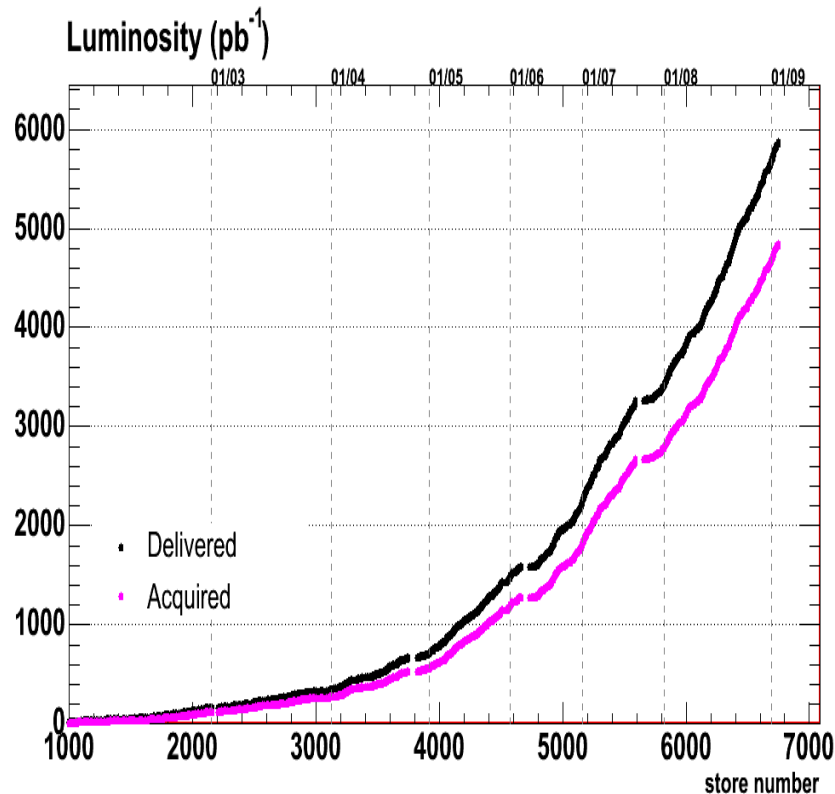
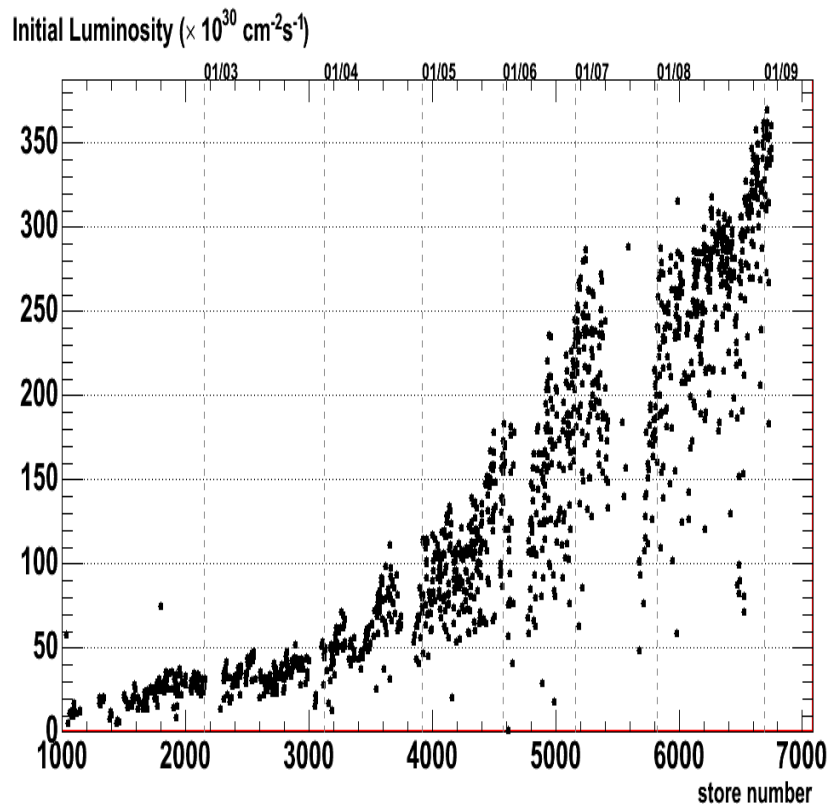
Rob Roser
For the CDF Collaboration

Thanks to Rick Snider and Donatella Lucchesi
for their help in preparing this talk.

Physics Analysis Strategy

- CDF Analysis is Conference Driven
- We work to have a "suite" of results ready for the winter "ski conferences" and the main Summer conference
- Just prior to those periods, new data is added to the analyses
- In between, analyzers are busy perfecting their craft, optimizing cuts and working to understand backgrounds and systematics better
- We have to have a computing system that can function with this type of "peakish" demand

The Accelerator Complex is Running Well!



Total Delivered 5.8 fb⁻¹

Total Acquired 4.8 fb⁻¹

Eff: 82.5% acquired and 77.4% good

Issues that Drive the Computing Model

- Computing demand
 - Raw data logging rate, total data volume
 - Complexity / sophistication of analysis
 - Number of people performing analysis / number of analyses
- Computing infrastructure and operations
 - Budget constraints
 - Evolving grid infrastructure, access policies
 - Access after LHC turn-on?
 - Number of people available for operations support
- In general, the computing problem becomes more difficult with time due to increasing demand and declining effort.
- Must evolve and adapt to meet these challenges.

Strategies to Deal With

- Manage demand via highly centralized, incremental data processing model
 - Allows most cost effective use of CPU
- Expand use of grid-based resources
 - Leverages effort used to create common tools
- Simplify systems, automate operational procedures
 - Reduce cost of systems and effort required to run them
- Increase uniformity of infrastructure
 - Both hardware and software

Computing Unlike Detector is not Static

- CDF constantly adapting
- Changes we have made in the past year....
 - Infrastructure
 - Consolidation of on-site CPU resources under Fermigrid
 - Retirement of aging hardware
 - Migration of data to higher density tape technology
 - Operations
 - Improved MC processing model: luminosity profile scaling
 - Saves factor of 5 in processing relative to run-based scaling
 - Calibration automation improvements
 - Eliminating manual steps required for each 6-12 week production cycle
 - Production processing improvements
 - Reducing time to get data to tape and recover from processing errors

Computing Demand Model

CPU demand model has two components

- Production activities
 - Reconstruction, data reduction, Monte Carlo simulation
 - Completely centralized / coordinated processing
 - Demand scales with data logging rate
- Analysis
 - Decentralized, largely uncoordinated activity
 - Demand scales with total data volume (at worst)
 - The number of people / analyses
 - Increasingly important with time!

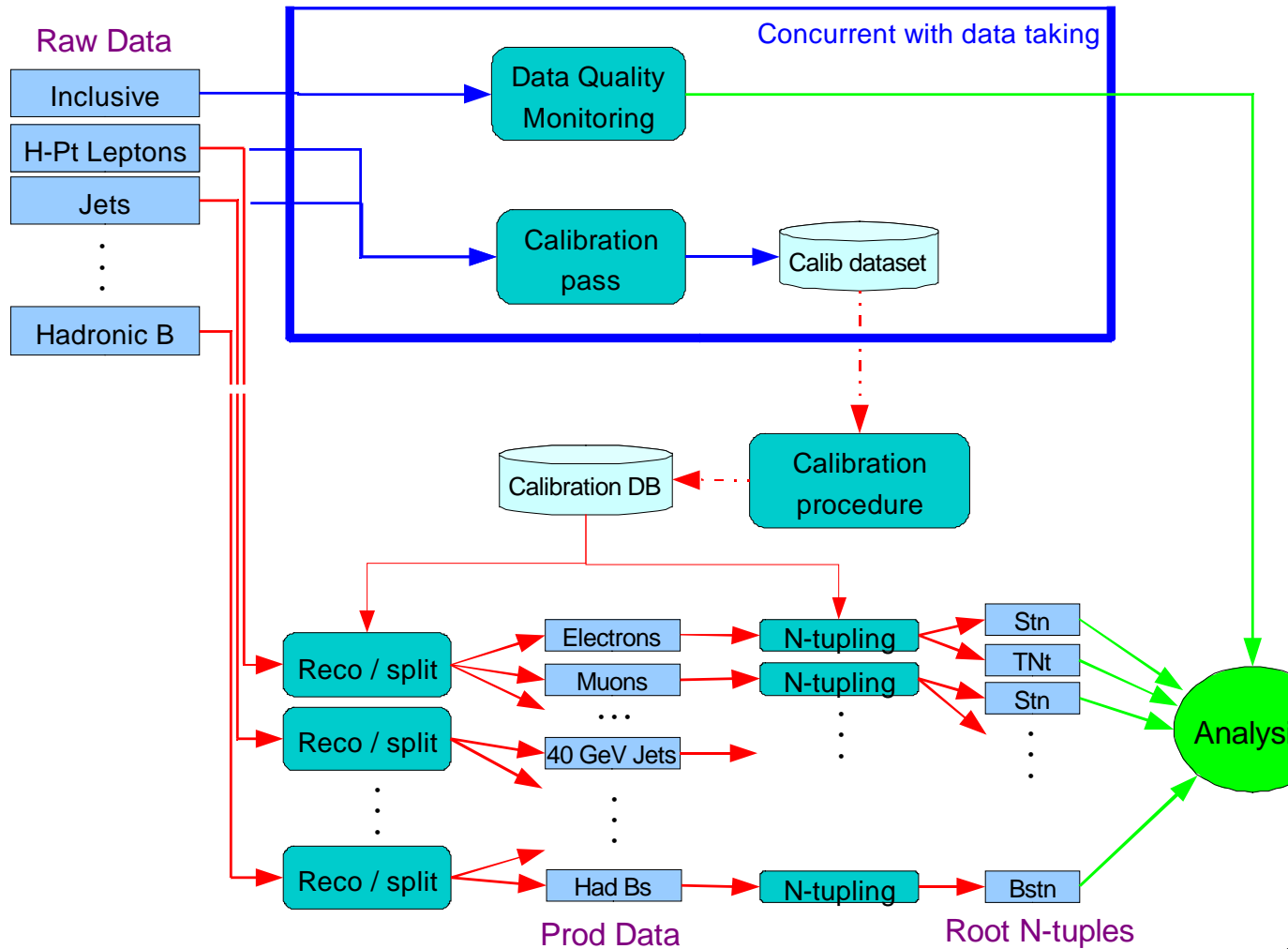
Raw data production model

- Goals of offline production operations
 - Deliver data required for analysis as close to data taking as possible
 - Final compressed datasets from reconstructed raw data
 - Ensure production is not the limitation in the rate of physics output
- The processing problem
 - Log data at rate of 5 - 7 M events/day
 - Event Complexity increases with increasing luminosity, but logging rate is \sim independent of accelerator luminosity
 - Calorimeters require re-calibration every \sim 3 months
 - Need to accumulate \sim 150+ M events to calibrate (though not all used for calib)

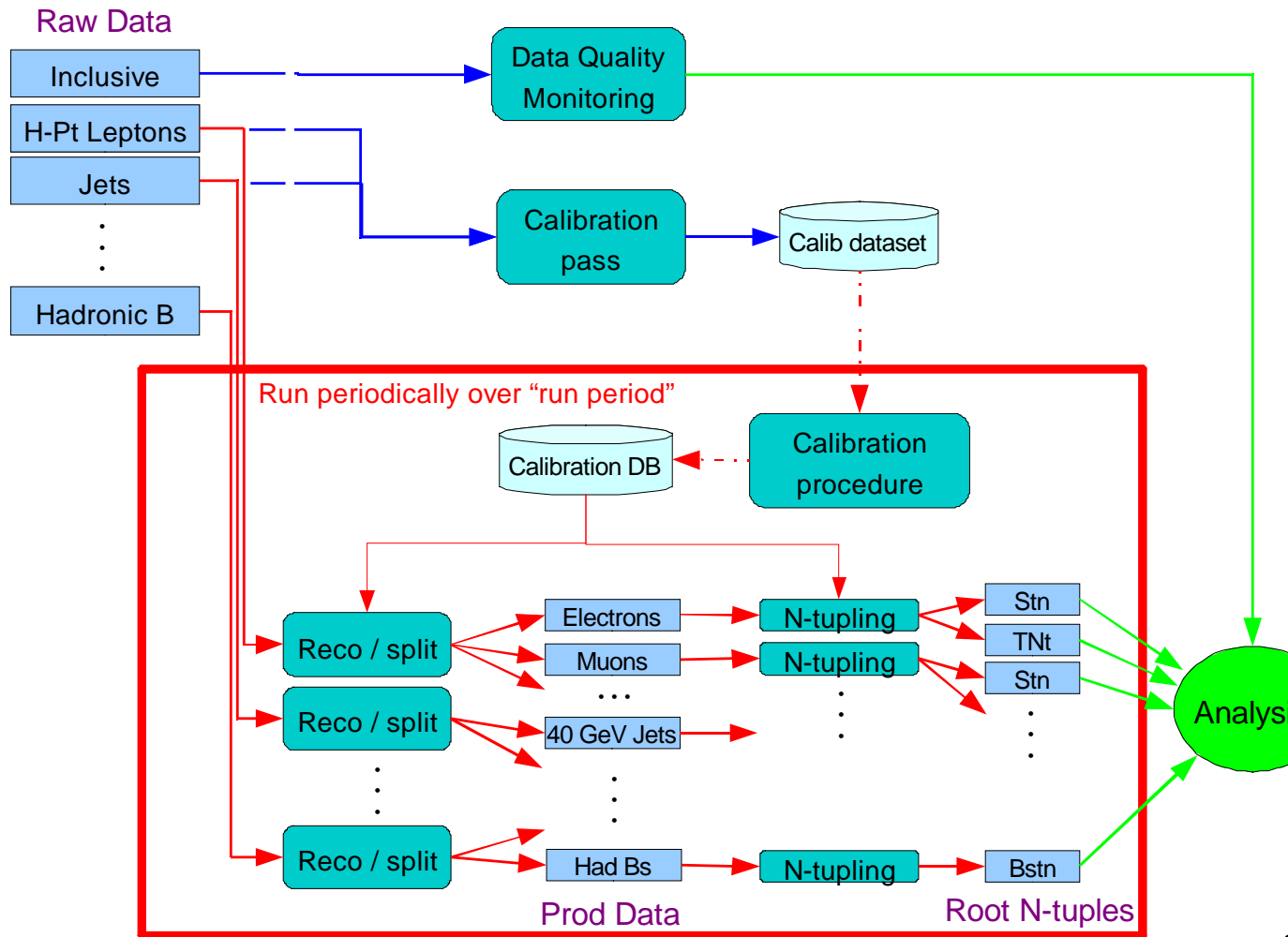
Raw data production model (2)

- Strategy
 - Divide data into "run periods" of 4 - 10 weeks
 - Typically 200 - 400 M events
 - Process data by run period
 - Calibration, raw data reconstruction, ntuple creation
 - Analyses use multiple run periods as needed

Raw data production model



Raw data production model



The production cycle

- Detector calibrations

- Process about 30% of the raw data within a few days following data taking
- Calculate calibrations and perform validation for each run period

Typically calibration completed 3 - 4 weeks after end of run period

- Raw data production

- Reconstruction of data
- Split data into datasets into physics datasets based upon triggers
 - 42 full + 9 compressed datasets

Typically completed 3 - 6 weeks after calibrations ready

The production cycle(2)

- Ntupling
 - This is the compressed and reconstructed data that all physics analyses are based upon
 - Performed on production output (after splitting)
 - Prioritize processing to do most important first
 - Three partially overlapping flavors: standard, top, Bs
- Typically 2 - 3 days behind raw data production

Raw data operations

- Event reconstruction
 - Average processing time
 - ~2 sec/event across all streams and luminosities
(varies greatly event type)
- Data processed on-site
 - Past run periods processed on 600 node farm dedicated to CDF
 - Also used for calibrations, N-tupling and analysis
 - Recently moved to processing on Fermigrid-based farms
 - Better optimizes CPU utilization
- Data re-processing
 - About 30% of data is processed twice as part of production cycle
 - Once for calibrations, once for physics datasets
 - The experiment has no plans for large scale re-processing

Monte Carlo data production

- The “old” MC data production model
 - Run-based MC that takes into account detector configuration and luminosity
 - Required continuous MC production operations coordinated with data taking
 - Deemed important early on given the changing beam and detector conditions
 - Final MC sample for an analysis could not be completed until the data was defined and complete
 - Big Conference periods were very stressful as we scrambled to get the data and MC ready so that analyzers could do their work.

Monte Carlo data production

- Changing the production model for new MC
- Need to change strategy with increasing data set size
- The new MC production model
 - Luminosity profile scaling
 - Generate MC asynchronously with data taking
 - Allows better scheduling of CPU usage
 - Significantly reduces amount of MC needed relative to run-based approach
 - Possible because both the detector configuration and the accelerator performance is very stable

Analysis Demand Model

- Separate analysis into several categories
 - "Core" analyses (as defined in the Tevatron Collider Experiment Task Force Report, Dec., 2005)
 - "Other" analyses
- Core analyses
 - Assume these are always fully staffed, so computing demand remains high
 - Some evolution in the analyses
 - More complex / sophisticated algorithms
 - Better procedures or more CPU efficient algorithms
 - All current data production activities needed to support core analyses
- Other analyses
 - Staffed with remaining effort
 - Demand scales with the number of people working on non-core analyses

Core Analyses

- As defined by the Tevatron task force report...
 - Identified "core" analyses that formed the basis of the justification for extended running of the Tevatron:
 - Measurement of Δm_s or limit on B_s mixing;
 - Measurement of $\Delta\Gamma_s/\Gamma_s$;
 - Limit on the branching ratio of the process $B_s \rightarrow \mu^+\mu^-$;
 - High precision measurement of the W boson mass;
 - High precision measurement of the top quark mass;
 - Measurement of single top production cross-section;
 - Search for the Higgs boson both in the Standard Model and SUSY scenarios;
 - Searches for SUSY in the "golden" mode Gaugino-neutralino with tri-leptons;
 - Searches for SUSY in the "golden" mode Squark-gluino with multijets plus missing transverse energy;
 - Searches for high mass resonances in the e^+e^- , $\mu^+\mu^-$, $\gamma\gamma$ and jet-jet invariant mass spectra (sensitive to Large Extra Dimensions, Z' and other processes not present in the Standard Model);

Analysis computing

- Computing requirements scale with:

- Full data set size
- Complexity of analyses
- Number of people / analyses

Computing problem becomes harder with time

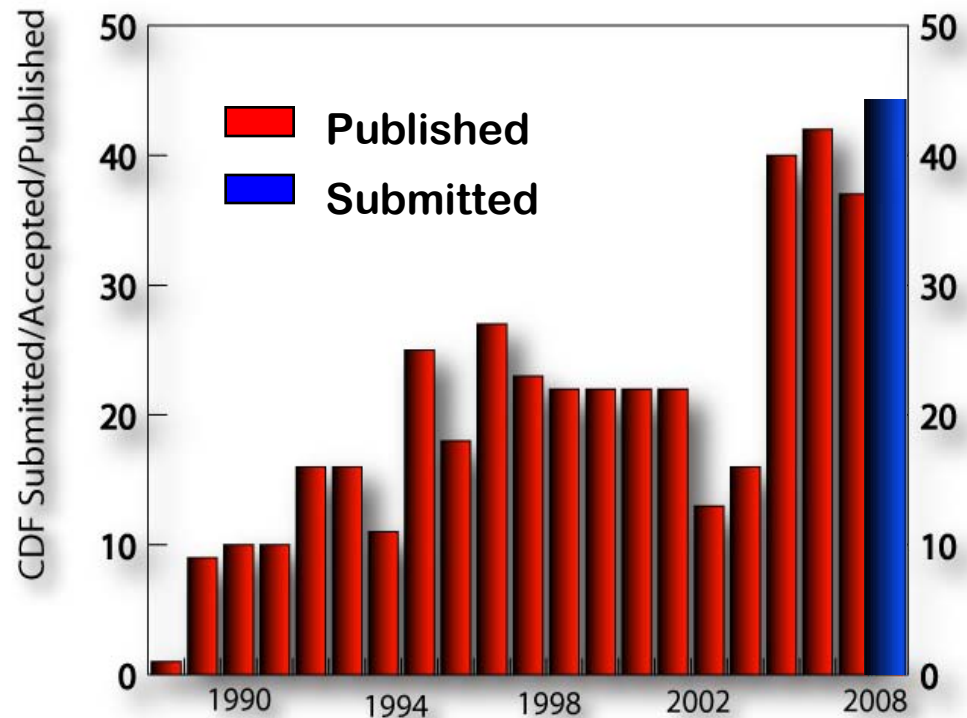
- Facilities

- 5k CPUs on-site for data intensive analysis
 - Shared with production activities
 - Some large datasets also located at INFN-CNAF
- Off-site computing also available for CPU intensive analysis
 - Matrix element analysis, pseudo-experiments, etc.

Period	Start	End	Lum (pb-1)	Events (M)	N-tuples ready
13	May 13, 07	Aug 4, 07	317	545	Nov 29, 07
14	Oct 28, 07	Dec 3, 07	45	59	Feb 21, 08
15	Dec 5, 07	Jan 27, 08	159	210	Apr 7, 08
16	Jan 27, 08	Feb 27, 08	142	168	May 21, 08
17	Feb 28, 08	Apr 16, 08	188	235	Jun 6, 08
18	Apr 18, 08	Jul 1, 08	407	436	Oct 25, 08

Analysis computing

- Is it all effective?
- The bottom line is the physics that CDF produces
 - Another 50+ new results at 2008 Summer conferences
 - 43 publications in 2008
 - 21 Submitted
 - Analyses keeping pace with the data sets!



Measured on-site demand

Production and Ntupling

- Expect to remain constant through end of data taking
- After that, depends on analysis needs

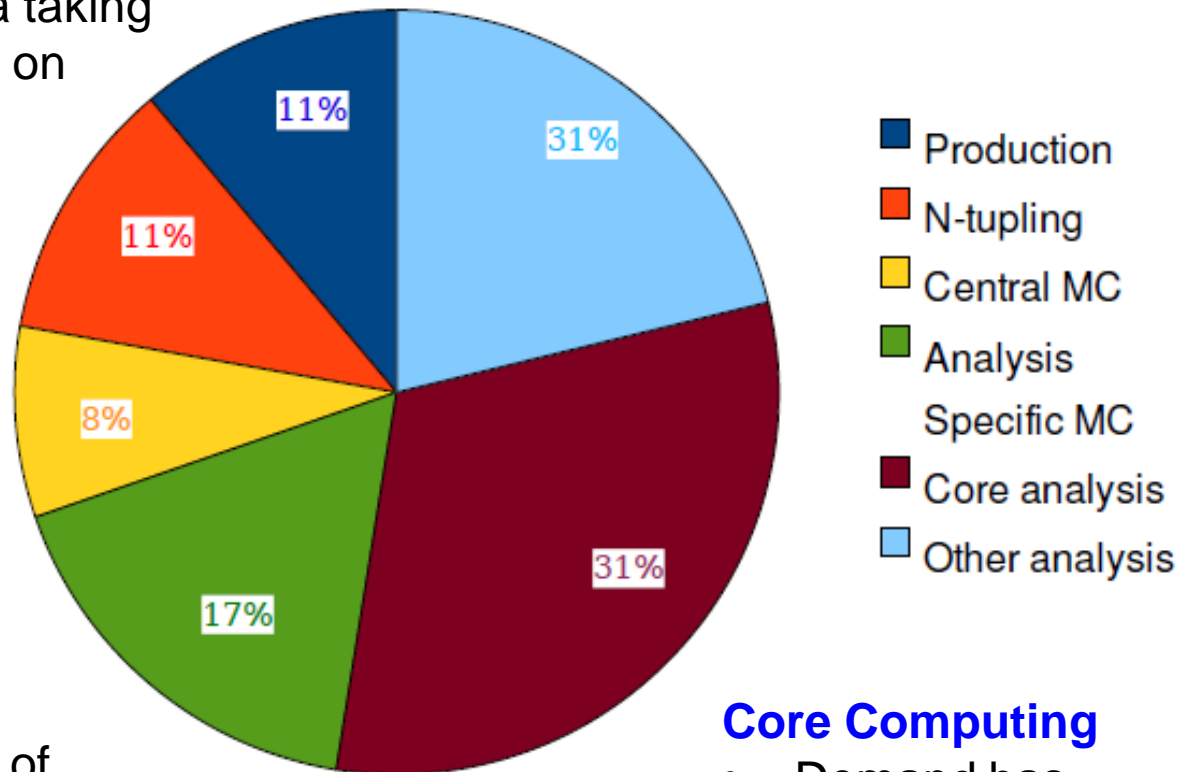
Core MC

- tied to size of data set

Non Core Computing

- Scales with number of people
- Will fall off as LHC ramps up

CDF On-site CPU usage

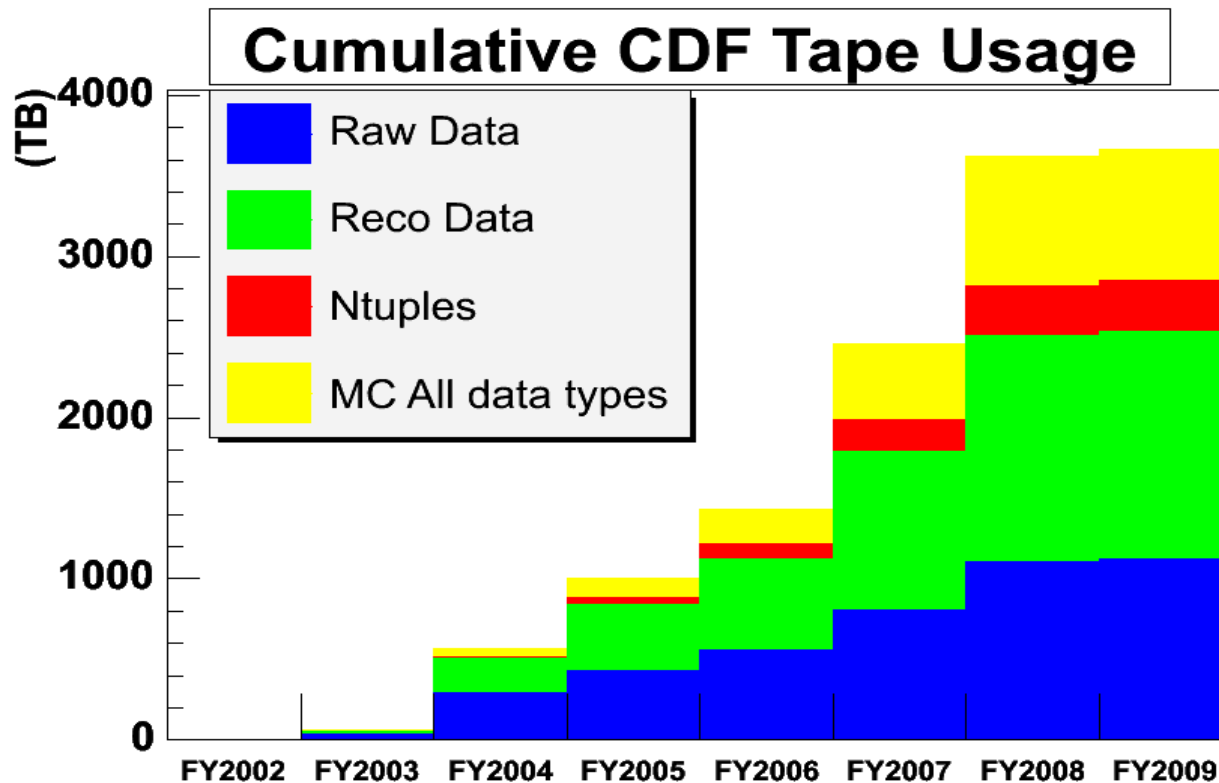


Core Computing

- Demand has remained constant
- Final results ~2 years after last data

Data volumes

- Data on tape
 - Total of 3.7 PB
 - Raw data
 - 7.9 billion events
- Monte Carlo data
 - 4.6 billion events
 - Includes a combination of centrally produced MC and analysis-specific MC



Summary

- The CDF offline is successfully meeting the physics needs of the experiment
 - Due to the hard work of many collaborators at Fermilab and around the world
 - A close and productive collaboration with the Computing Division has been critical to this success
- Will ensure continued success by working to improve the systems, increase efficiency and reduce the effort required to conduct computing operations.

Some Physics Highlights

- ✓ **Observation of Bs-mixing**
 - $\Delta m_s = 17.77 \pm 0.10$ (stat) ± 0.07 (sys)
- ✓ **Observation of new baryon states**
 - Σ_b and Ξ_b
- ✓ **WZ discovery (6-sigma)**
 - Measured cross section 5.0 (1.7) pb
- ✓ **ZZ observation**
 - 4.4-sigma
- ✓ **Single top evidence (5-sigma sens./2.2 fb⁻¹)**
 - cross section = 2.2 (0.7) pb
 - $|V_{tb}| = 0.88 \pm 0.14$ (exp.) ± 0.07 (th.)
- ✓ **Observation of new charmless B_s → hh states**
- ✓ **Observation of exclusive/diffractive production**
 - Di-jets, W/Z, charmonium, etc
- ✓ **Observation of D⁰-D⁰bar mixing**
- ✓ **Measurement of Sin(2β_s)**
- ✓ ...
- ✓ **Precision W mass measurement**
 - $M_{W_cdf} = 80.413$ GeV (48 MeV)
- ✓ **Precision Top mass measurement**
 - $M_{top_cdf} = 172.4$ (1.6) GeV
- ✓ **W-width measurement**
 - 2.032 (.071) GeV
- ✓ **Extended exclusions on BSM**
- ✓ **Continued improvement in Higgs Sensitivity**
 - Exclusion of 170 GeV Higgs

MC data production operations

- MC data generated
 - 1.1 G events produced last year
 - Some periods of concentrated production during "MC attacks"

