

CMS Experiences In CSA06

A Computing, Software & Analysis
Challenge in 2006

Jorge L. Rodriguez
University of Florida
Reporting for the CSA06 Group

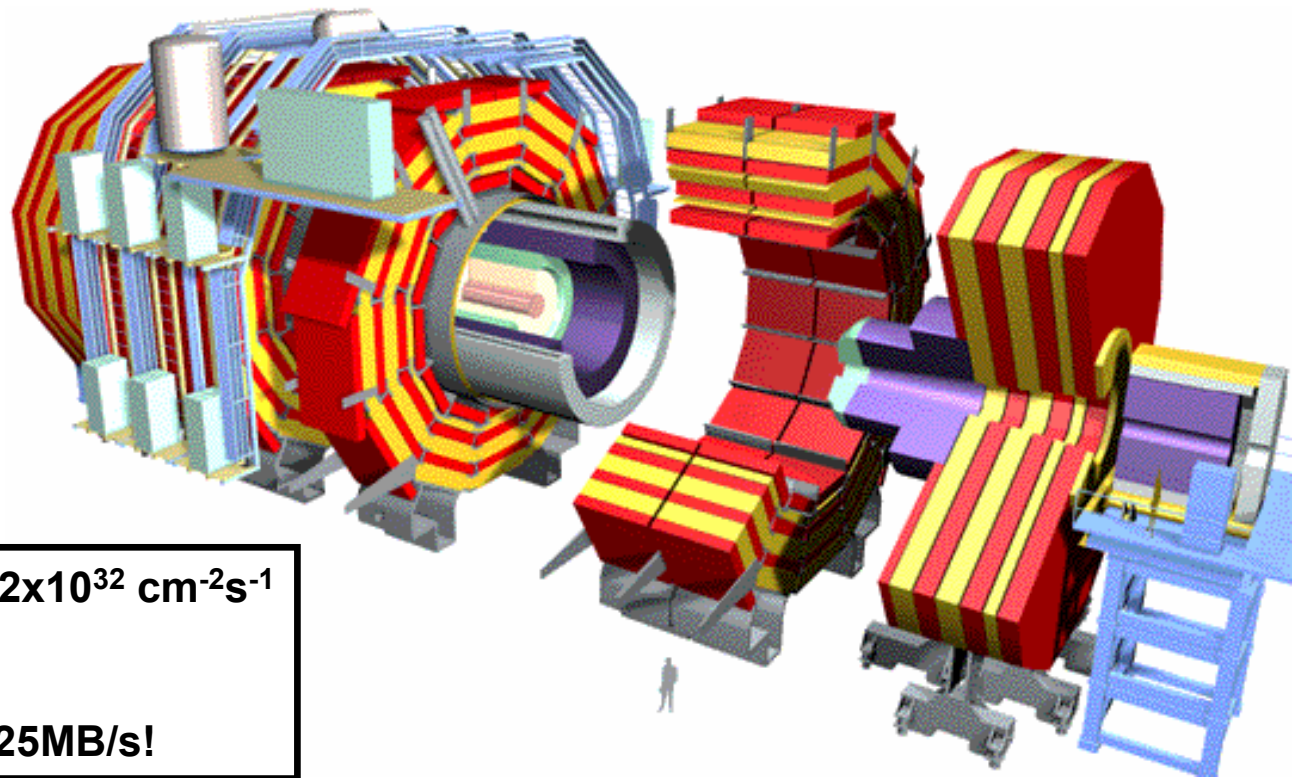




The CMS Experiment

CMS “Compact Muon Solenoid” 1 of 4 experiments at the LHC

- Millions of channels + High Luminosity + processing ~ 10 PB/yr
- Lots of experimenters more than 2500 located all over the world



At startup luminosity of $2 \times 10^{32} \text{ cm}^{-2}\text{s}^{-1}$
Trigger rate ~ 150 Hz
Events size ~ 1.5 MB
Data written to tape @ 225 MB/s!



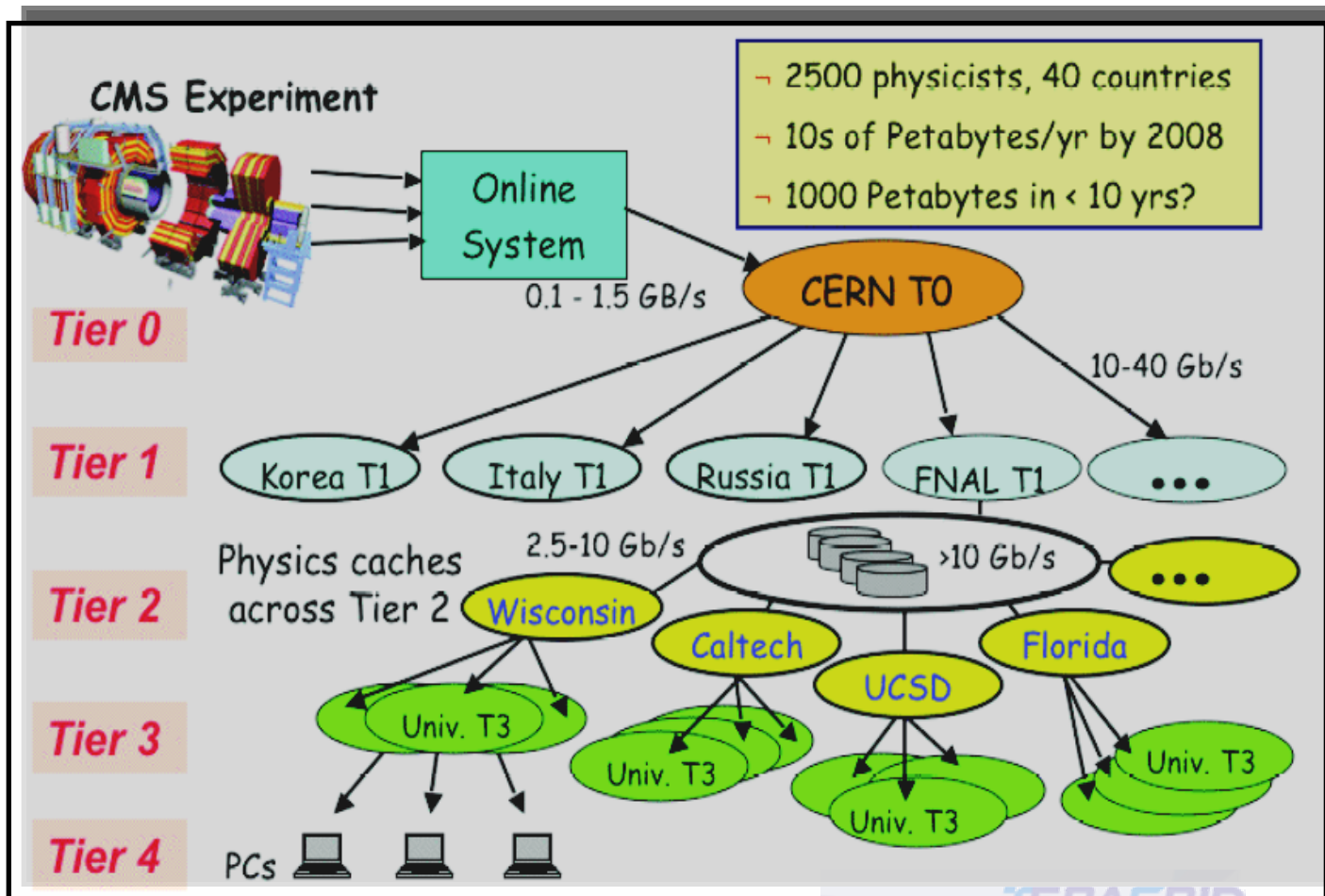
CMS' Distributed Computing Model

WLCG



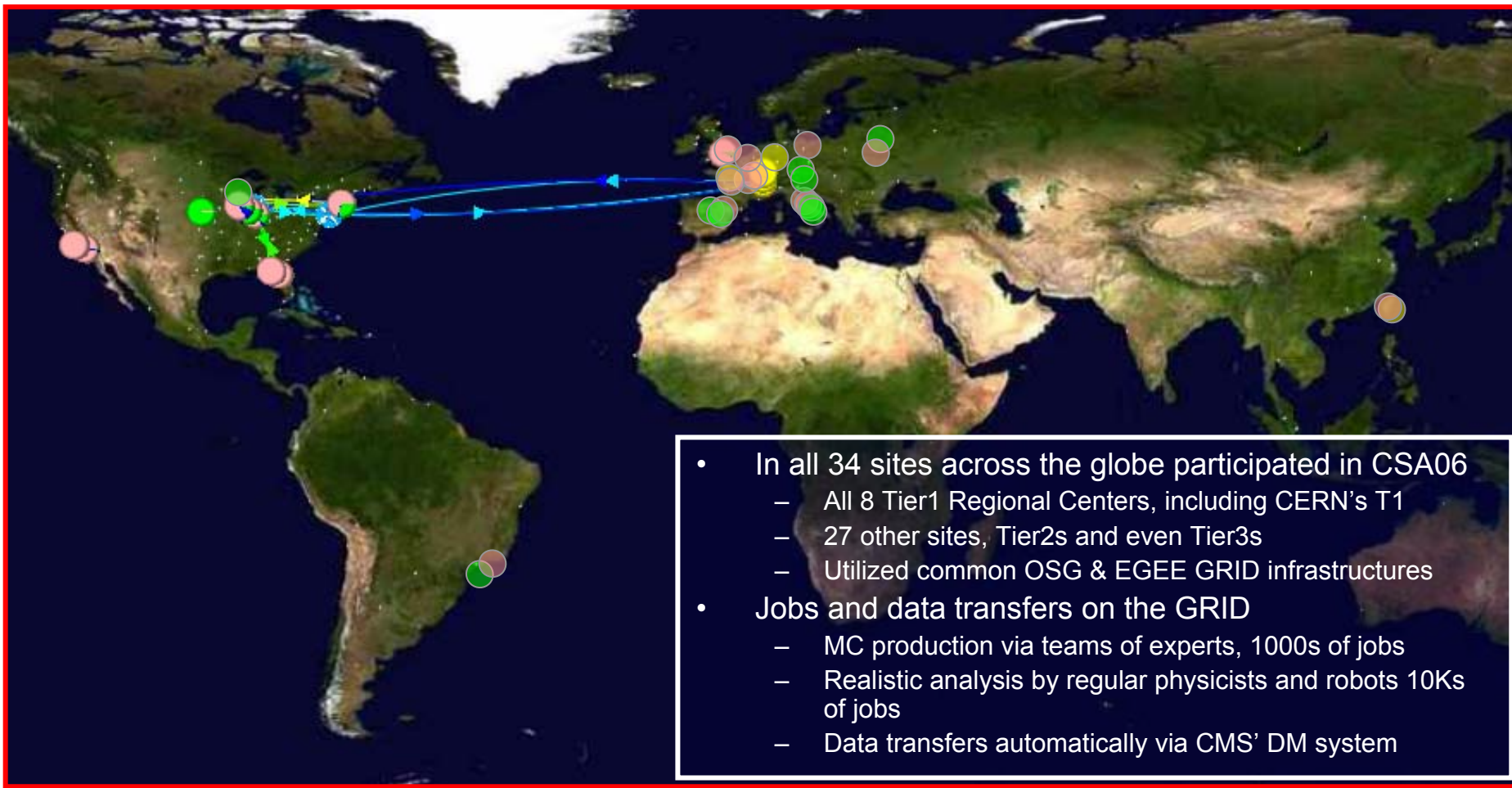
eGEE

Enabling Grids for E-science in Europe





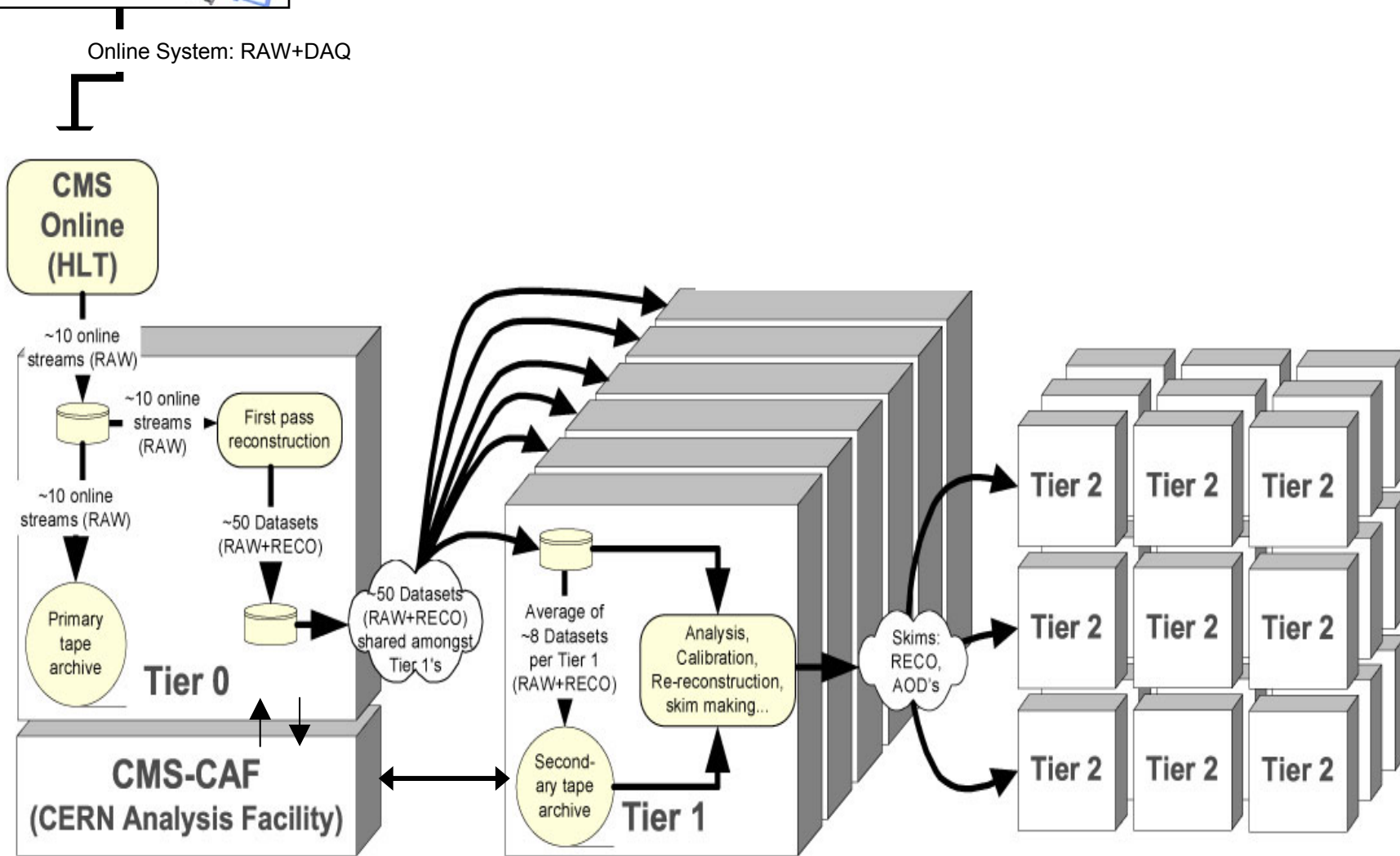
CMS Computing on the GRID



- In all 34 sites across the globe participated in CSA06
 - All 8 Tier1 Regional Centers, including CERN's T1
 - 27 other sites, Tier2s and even Tier3s
 - Utilized common OSG & EGEE GRID infrastructures
- Jobs and data transfers on the GRID
 - MC production via teams of experts, 1000s of jobs
 - Realistic analysis by regular physicists and robots 10Ks of jobs
 - Data transfers automatically via CMS' DM system



CMS Data Flow





CMS Event Data Model

Evt Format	Data Contents	Evt Size (MB)	Evts/yr	Volume/yr (PB)
DAQ+RAW	Detector data + L1 trig. results	1-1.5	1.5×10^9	-
RAW	Detector data after online formatting + L1 trig results + HLT trig. Bits	1.5	3.3×10^9 2 copies+overlap	5.0
RECO	Reconstructed objects, tracks, jets... + all hits & clusters	0.25	8.3×10^9 2 copies + 3 reprocessing	2.1
AOD	Reconstructed objects, tracks, jets... + hits & clusters + small quantities of localized hit info	0.05	53×10^9 4 versions + 8 copies at Tier1	2.6
TAG	Run/evt number, high level physics objects used to index events	0.01	-	-
FEVT	Term used to denote RAW+RECO	-	-	-



The CSA06 Challenge: What is it?

- **A 50 million event exercise to test the workflow and dataflow associated with the data handling and data access model of CMS**
- **A 25% capacity test of what we will need in 2008**



Overall Goals of CSA06

- Demonstrate designed workflow and dataflow
- Demonstrate Computing-Software synchronization
 - Go smoothly through one or more CMS Software updates
- Demonstrate production-grade reconstruction software
 - Includes calibration; detector performance
- Demonstrate all cross-project actions
 - Determination and use of calibration/alignment constants including insertion and extraction and offline use of said constants via distributed constants database system
 - The HLT exercise: Split pre-challenge samples into multiple “tagged” streams and process these through complete DM system
- Provide services to a wide user community
 - Not just robotic GRID submissions
 - Support local and remote GRID users



Quantitative Goals and Metrics

- Site participation
 - Tier1s: Goal is all, but more than 5
 - Tier2s: Goal is 20 but more than 15
- Tier0 Processing farm
 - Number of weeks of sustained running: Goal is four weeks
 - Tier0 efficiency: Goal of 80% but more than 30%
- Data Management and Movement System
 - Data transfer rates from Tier0 to Tier1 to tape per site: Goal more than 50% of site specific capacity
 - Data transfer rates from Tier1 to Tier2: Goal is 20 but more than 5 MB/sec/site
- Physics Analysis Jobs on the GRID
 - Running jobs at Tier1 and Tier2 [2hr jobs/day]: Goal is 50K/day but more than 30K
 - GRID job efficiency: Goal 90% but more than 70%



Binary Goals and Exercises

- Are global systems operational?
 - Does the Data Mgmt system work? Focus on T0 → T1 → T2 transfers
 - Does Constants DB system work? Can we read DB offline & remotely can we insert new constants into system?
 - Can we run analysis & skim jobs via CRAB on the GRID?
- Exercises designed to test Workflows in CMS
 - Alignment Exercises
 - Get corrections from misaligned datasets (Tracker, & Muon Sys)
 - Inject corrections into DB and apply in Re-reconstruction
 - Calibration Exercises
 - ECAL, HCAL misalignment and re-reconstruction as above
 - ECAL, HCAL: ϕ symmetry calibration, HCAL jet corr. functions...
 - Analysis Exercise Demonstations
 - Extraction of signals: dimuons distributions, Higgs mass plots...
 - T-Tbar, SUSY/BSM, Standard model background and other studies

Defining the CSA06 Challenge

66M MC events simulated up to detector digitization with no-pileup

HLT tagged events were based on MC truth tables

Tier1&Tier2s process GRID analysis jobs submitted either by job robots or by CMS physicists using CRAB (CMS' GRID job builder/submitter)

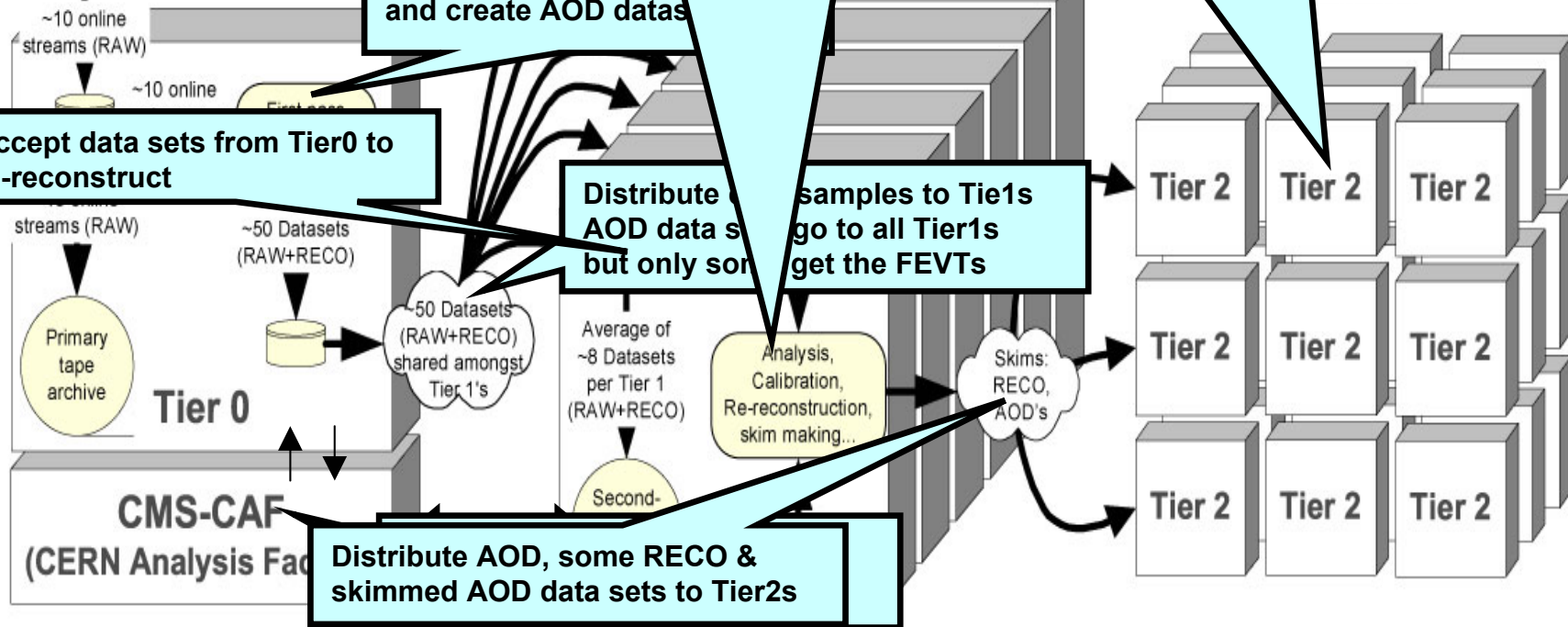
Run re-reconstruction with updated constants and CMS software. Also run skim jobs on AOD data sets and create AOD data sets

Accept selected data sets from Tier1s and store in local SE

Accept data sets from Tier0 to re-reconstruct

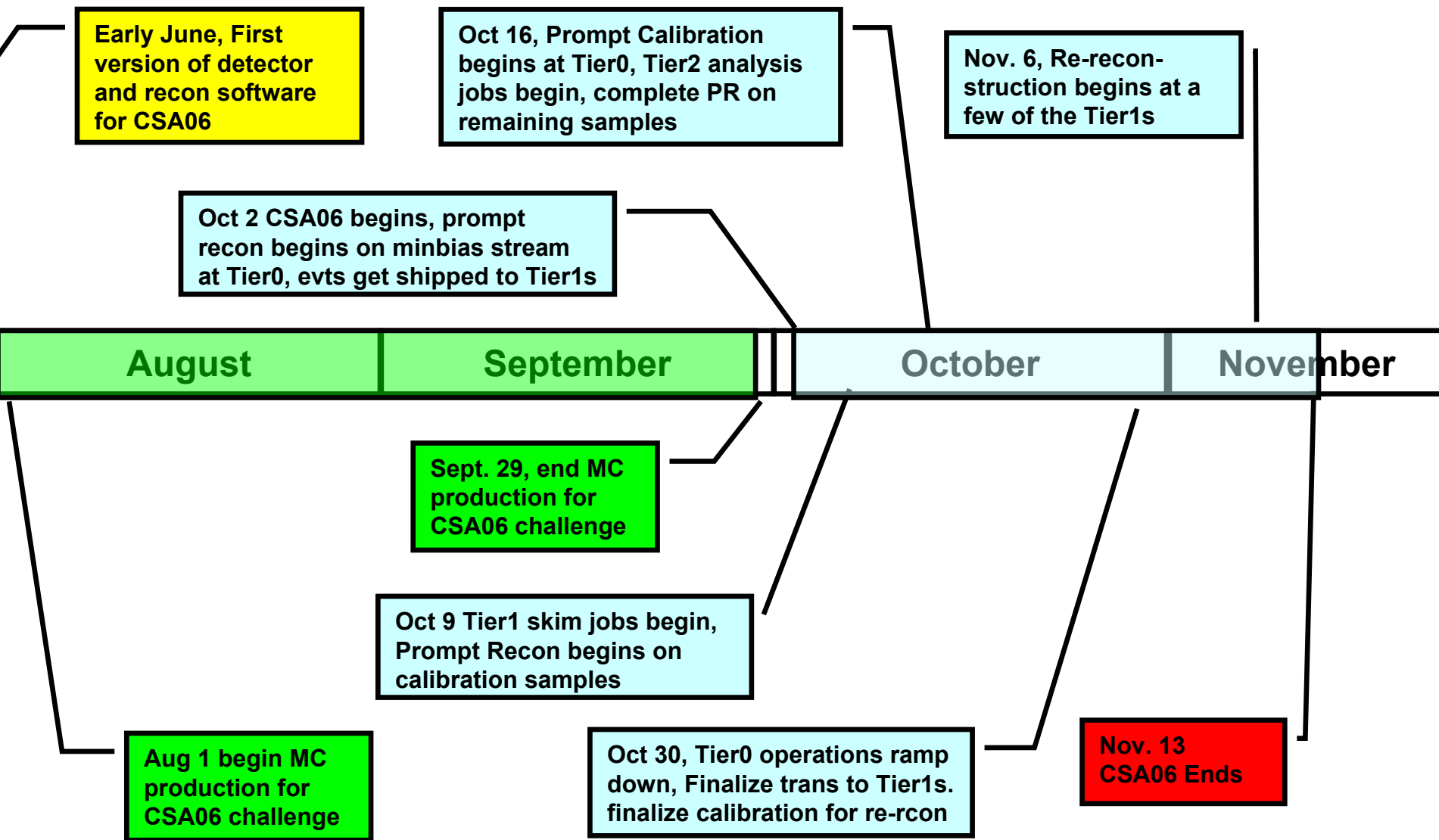
Distribute AOD data sets to all Tier1s but only some get the FEVTs

Distribute AOD, some RECO & skimmed AOD data sets to Tier2s





CSA06 Time Line





Status as of 10/31/06



CSA06 Status Report

“Many pieces of CSA06 already successful”

Michel Ernst @ RRB Meeting on 10/24/06

- Some Highlights

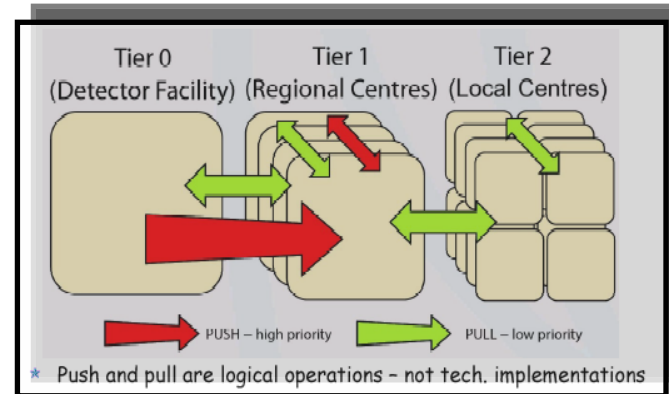
- 66M events simulated with new CMSSW (Aug-Sept)
- 207M events reconstructed at Tier0 with 100% uptime for 4 weeks
- All Tier0 workflows completed successfully
 - <50 Hz> for first 3 weeks, increased to 100Hz the last 5 days
 - New CMS software stable and performs well
 - Calibration output produced & registered into DMS quickly & efficiently
- More than 900TB of data transferred to Tier1 & Tier2s via the DMS
- Skim & Analysis jobs underway & running very well
 - Average about 20,000 jobs/day across the GRID
 - Physics Analysis jobs success rates @Tier2s are excellent



CMS Data Management System

A Distributed Data Management System

- Identification & location of data
- Worldwide movement of data replicas
 - Tier0 ↔ Tier1 ↔ Tier2 ↔ Tier0...
 - Volume ~ PBs ≈ O(10 M) files
 - Transfer speeds ~ 5 Gb/s
- Data collections
 - Datasets: collections of data blocks
 - Blocks: collections of data files ~ 10TBs
 - Files: collections of events ~ GBs
- Components include
 - PhEDEx:
 - Hart of the system
 - Manages data at the block level
 - DLS: Data location Service
 - DBS: Data Bookkeeping System



PhEDEx gives user scontrols over movement of replicas via a user interface. It also provides monitoring of transfers and of volumes.

The screenshot shows the PhEDEx Production Status web interface. It displays a table of transfer rates for the last 7 days. The table includes columns for To, From, Files, Total, Rate, Errors Expired, Avg. Est., and Avg. Est. The data is summarized in the following table:

To	From	Files	Total	Rate	Errors Expired	Avg. Est.	Avg. Est.
TL_RWTH_buffer	TL_DEST_buffer	2911	2.8 TB	4.3 MB/s	1934	12373	4.5 MB/s
TL_MT_buffer	TL_FINAL_buffer	2363	2.0 TB	3.4 MB/s	4418	4377	5.9 MB/s
TL_FINAL_buffer	TL_FINAL_buffer	439	362.4 GB	1.5 MB/s	903	4895	4.4 MB/s
TL_Spwn_JFCA	TL_FJC_buffer	1827	161.6 GB	1.3 MB/s	12969	15479	372.3 MB/s
TL_DEST_buffer	TL_CERN_buffer	1599	263.2 GB	454.4 MB/s	176	471	797.6 MB/s
TL_DEST_buffer	TL_FINAL_buffer	140	239.2 GB	414.7 MB/s	560	214	917.6 MB/s
TL_DEST_buffer	TL_FJC_buffer	42	123.5 GB	214.7 MB/s	52	-	9.3 MB/s
TL_Spwn_buffer	TL_FJC_buffer	151	113.3 GB	194.7 MB/s	1072	104	38.7 MB/s
TL_DEST_buffer	TL_RWTH_buffer	15	111.7 GB	103.8 MB/s	506	2	1.8 MB/s
TL_DEST_buffer	TL_ASSOC_buffer	18	57.0 GB	98.8 MB/s	-	-	4.4 MB/s
TL_DEST_buffer	TL_KAL_buffer	16	18.7 GB	34.1 MB/s	19	3	215.3 MB/s
TL_DEST_HSS	TL_DEST_buffer	72	15.5 GB	24.9 MB/s	-	-	103.3 MB/s
Total		8858	7.0 TB	12.1 MB/s	24888	41128	-/5

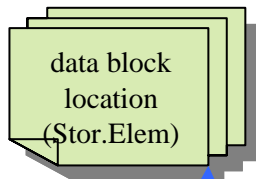


CMS Data Management System

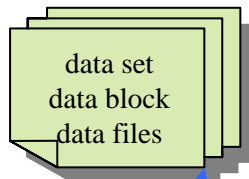
I. Fisk



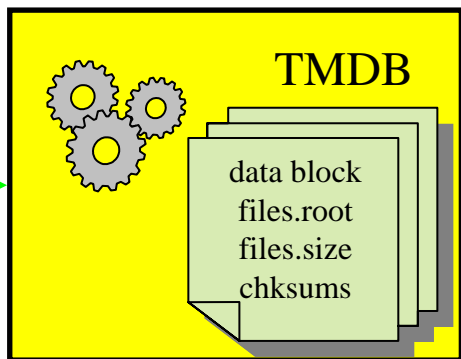
DLS



DBS



PhEEx.cern



FNAL

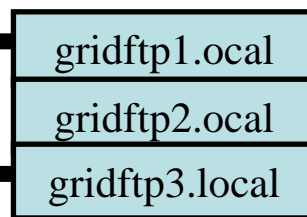
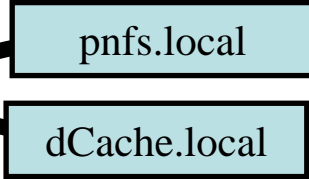
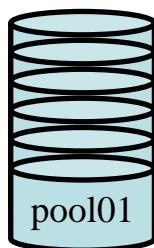
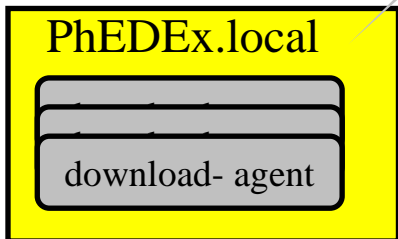
PhEEx.local

srmCache

`srmcp srm://<fnal> srm://<uftier2>`



UFlorida Tier2





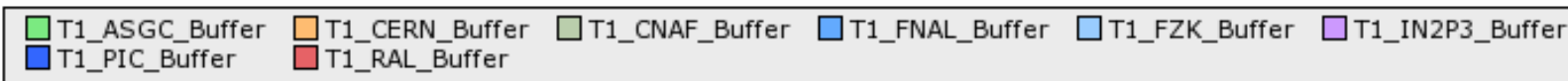
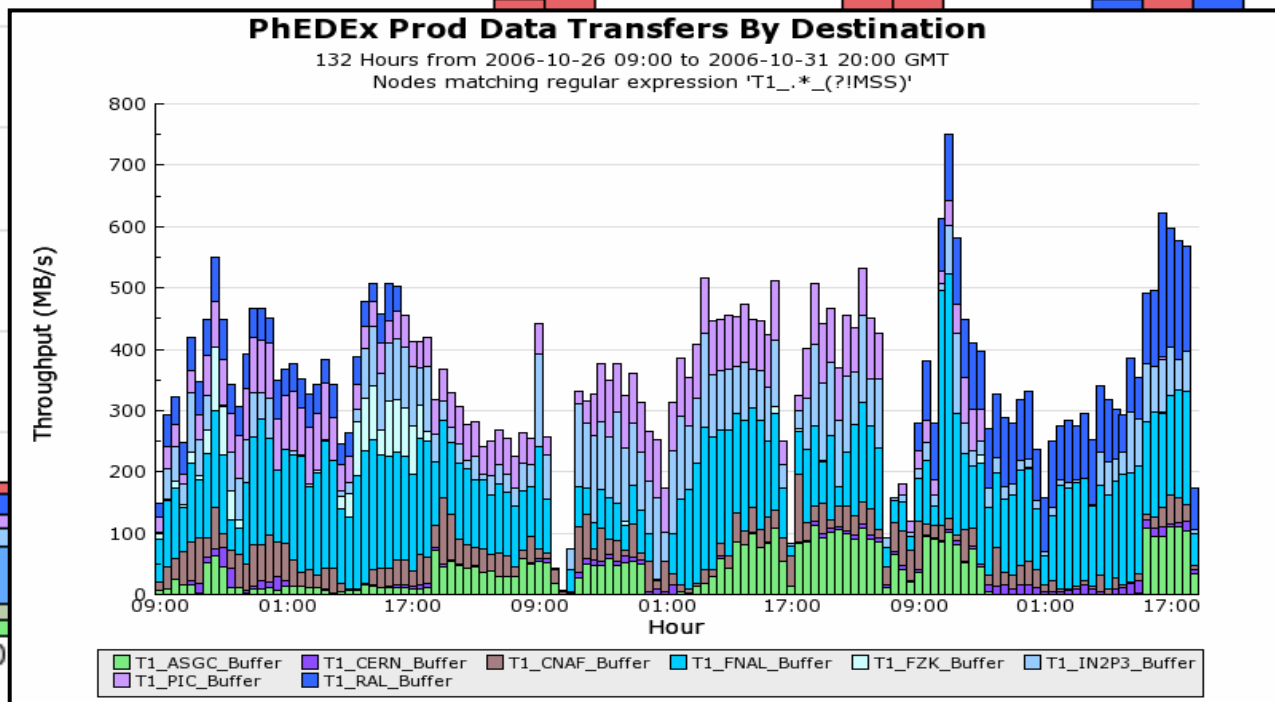
PhEDEx Tier1 Transfer Rate

PhEDEx Prod Data Transfers By Destination

30 Days from 2006-10-02 to 2006-10-31 GMT
Nodes matching regular expression 'T1_.*_(?!MSS)'

- All 8 Tier1 centers are represented registering transfers daily
- Problems with Storage Elements at some sites recovered quickly (< 18hrs)
- Over the last 4 weeks we've averaged ~200MB/s <day>
- Highest <hr> rate exceeded 700MB/s after on Monday!

Throughput (MB/s)



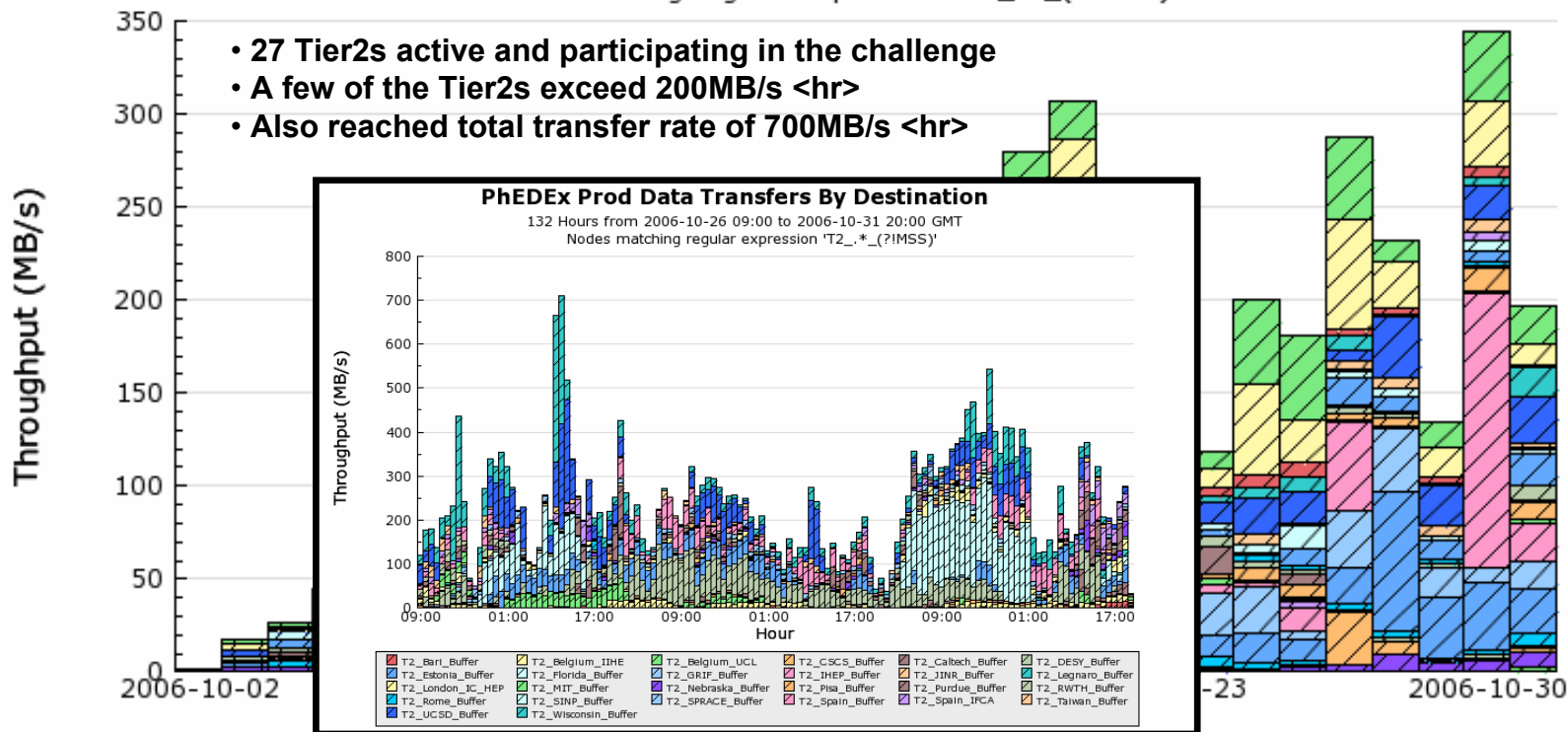


PhEDEx Tier2 Transfer Rate

PhEDEx Prod Data Transfers By Destination

30 Days from 2006-10-02 to 2006-10-31 GMT
Nodes matching regular expression 'T2_.*_(?!MSS)'

- 27 Tier2s active and participating in the challenge
- A few of the Tier2s exceed 200MB/s <hr>
- Also reached total transfer rate of 700MB/s <hr>



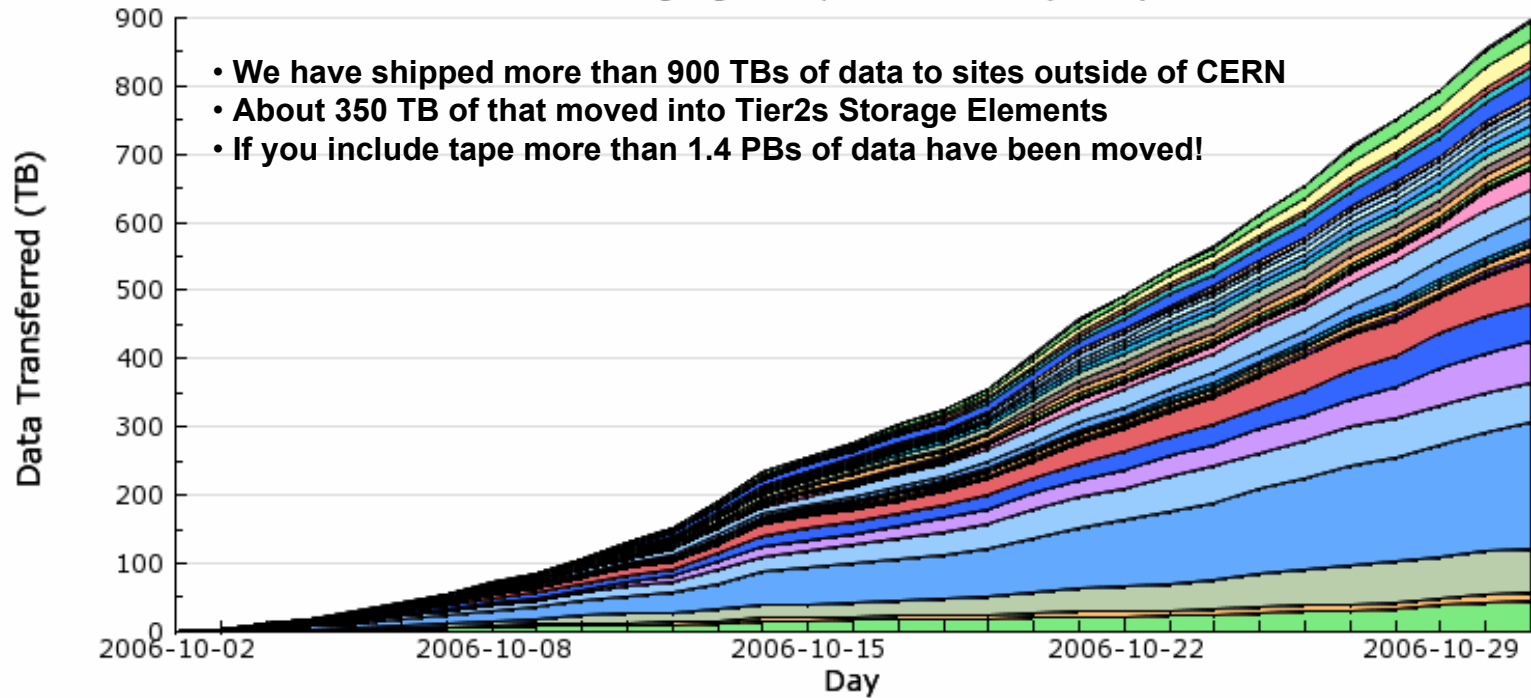
T2_Bari_Buffer	T2_Belgium_IIHE	T2_Belgium_UCL	T2_CSCS_Buffer	T2_Caltech_Buffer	T2_DESY_Buffer
T2_Estonia_Buffer	T2_Florida_Buffer	T2_GRIF_Buffer	T2_IHEP_Buffer	T2_ITEP_Buffer	T2_JINR_Buffer
T2_KNU_Buffer	T2_Legnaro_Buffer	T2_London_IC_HEP	T2_MIT_Buffer	T2_Nebraska_Buffer	T2_Pisa_Buffer
T2_Purdue_Buffer	T2_RWTH_Buffer	T2_Rome_Buffer	T2_SINP_Buffer	T2_SPRACE_Buffer	T2_Spain_Buffer
T2_Spain_IFCA	T2_Taiwan_Buffer	T2_UCSD_Buffer	T2_Wisconsin_Buffer		



Total Data Volumes Moved

PhEDEx Prod Data Transfers By Destination

30 Days from 2006-10-02 to 2006-10-31 GMT
Nodes matching regular expression 'T?_*_(?!MSS)'



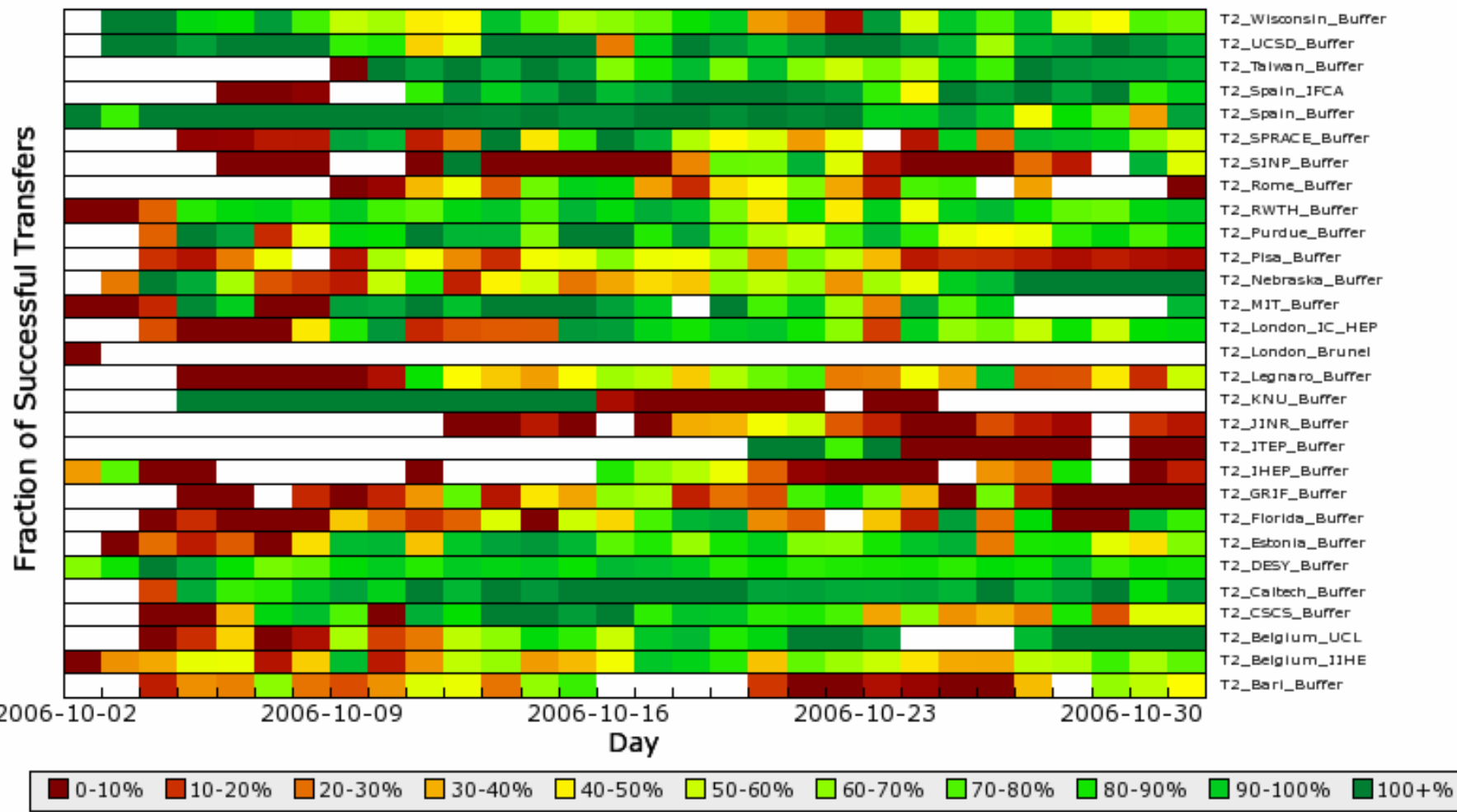
T1_ASGC_Buffer	T1_CERN_Buffer	T1_CNAF_Buffer	T1_FNAL_Buffer	T1_FZK_Buffer	T1_IN2P3_Buffer
T1_PIC_Buffer	T1_RAL_Buffer	T2_Bari_Buffer	T2_Belgium_IIHE	T2_Belgium_UCL	T2_CSCS_Buffer
T2_Caltech_Buffer	T2_DESY_Buffer	T2_Estonia_Buffer	T2_Florida_Buffer	T2_GRIF_Buffer	T2_IHEP_Buffer
T2_ITEP_Buffer	T2_JINR_Buffer	T2_KNU_Buffer	T2_Legnaro_Buffer	T2_London_IC_HEP	T2_MIT_Buffer
T2_Nebraska_Buffer	T2_Pisa_Buffer	T2_Purdue_Buffer	T2_RWTH_Buffer	T2_Rome_Buffer	T2_SINP_Buffer
T2_SPRACE_Buffer	T2_Spain_Buffer	T2_Spain_IFCA	T2_Taiwan_Buffer	T2_UCSD_Buffer	T2_Wisconsin_Buffer
T3_Minnesota_Buffer					



Operational Performance

PhEDEx Prod Transfer Quality By Destination

30 Days from 2006-10-02 to 2006-10-31 GMT
Nodes matching regular expression 'T2_.*_(?!MSS)'



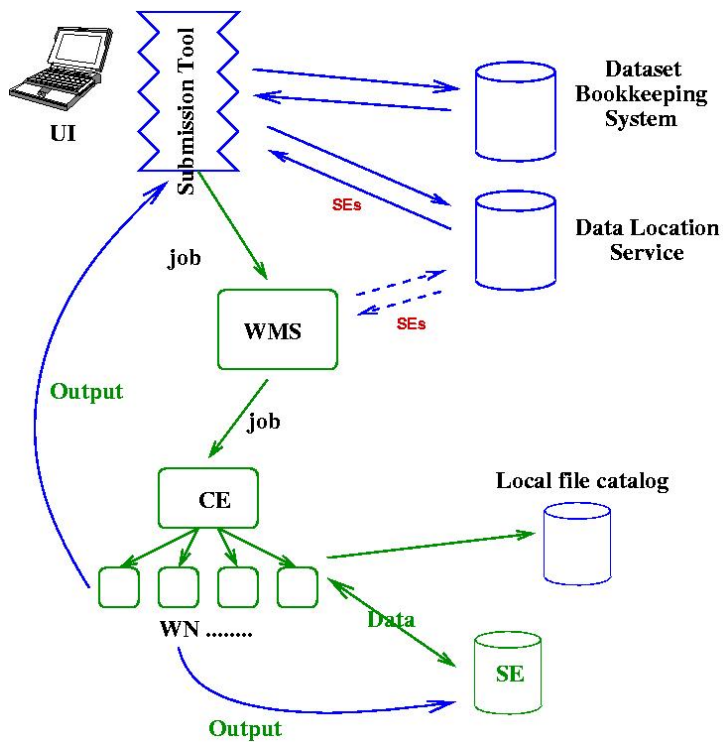


CRAB: CMS Remote Analysis Builder

An application to enable CMS analyses on the GRID

– User provides CRAB with:

- Dataset name, number of events
- Analysis code and job parameters: output file name & location etc...



– CRAB provides users with:

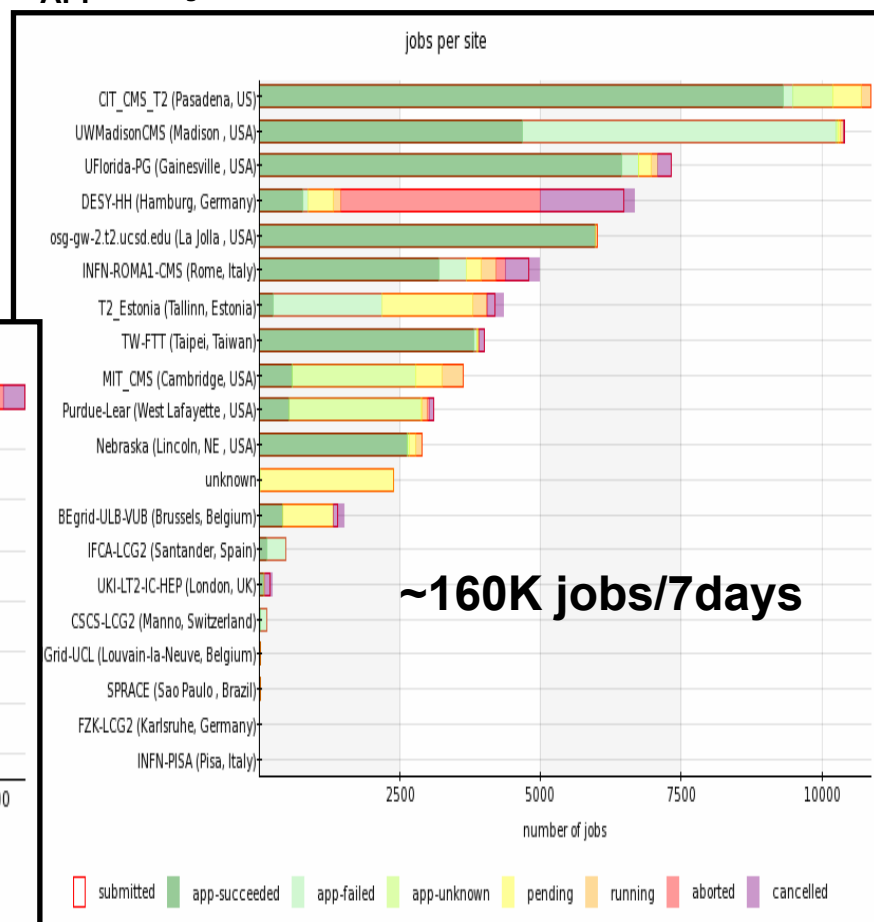
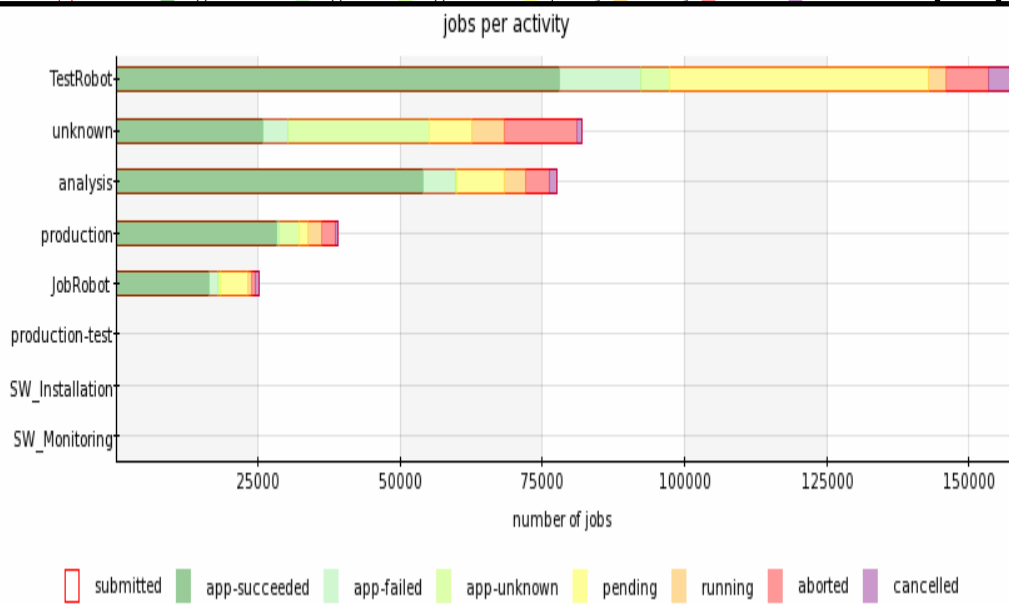
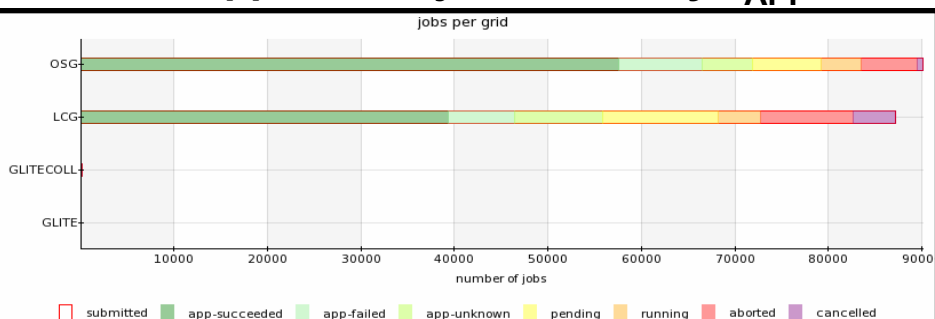
- Job creation and submission to GRID
- Data location and discovery
- Packaging of user code, bin, libs, data...
- Job control, scripts, shell wrappers
- Job monitoring and output management



Preliminary Look @ CMS GRID jobs

10/24/2006 – 10/31/06

- Overall GRID job efficiency $\epsilon_{\text{GRID}} = 85\%$: ϵ_{GRID} are jobs that get to site's worker node
- Overall application job efficiency $\epsilon_{\text{APP}} = 86\%$: ϵ_{APP} are jobs that return exit code 0





Status of Exercises

<https://uimon.cern.ch/twiki/bin/view/CMS/CSA06>

CSA06 < CMS < Twiki - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

https://uimon.cern.ch/twiki/bin/view/CMS/CSA06

Jump Search

[Edit](#)
[WYSIWYG](#)
[Attach](#)
[Printable](#)
[PDF](#)

You are here: [Twiki](#) > [CMS Web](#) > [CPTWikiHome](#) > CSA06

r186 - 01 Nov 2006 - 00:09:19 - Main.fisk

CSA06: Computing, Software, and Analysis challenge

Tier-0 Completes Operations, Analyses Well Underway

PHEX Prod Data Transfers By Destination
30 Days from 2006-09-29 to 2006-10-29 GMT

Jobs per minute
Oct 31 2006 09:25 UTC

Job availability

1PB transferred! (xfers last 48 h)

Frontier access up to end of CSA processing

Jobs in last 5 days

Wmunu skim
(Biasotto/Margoni/Torassa)

tau validation (Gennai)

Zmumu (Garcia)

Zee elec effic (Meridiani)

Misalignment (De Filippis)

Done

uimon.cern.ch



Summary

CSA06 already is quite a success story

- We have met or exceeded most of our targets and goals
 - Greater than 50 Hz processing at Tier0 (100% uptime for 4 wks)
 - Impressive transfer rates to Tier1s & Tier2s
 - All workflows and dataflows successful
 - CMS Data Management system shown to perform very well in realistic environment
 - Lots of data moved around the world, automatically
 - Limited only by data availability at the Tier0
- CSA06 ends in two weeks
 - Tier0 processing now completed (200M+ events processed)
 - More than 1.4 PBs of data moved by CMS' DMS
 - Analysis jobs are now running at full bore

THE END

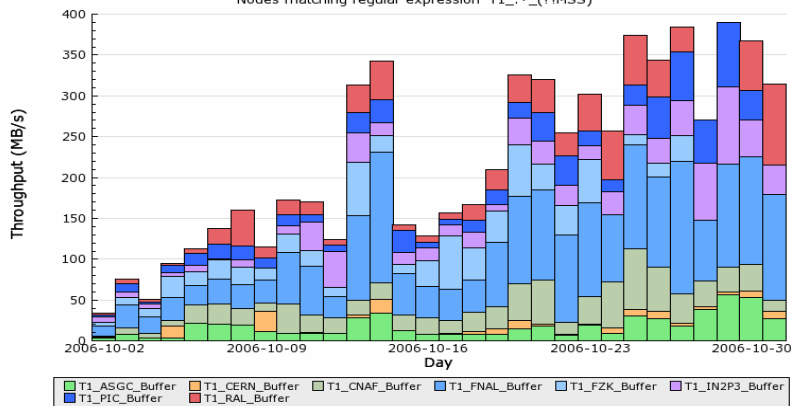
EXTRA SLIDES



PhEDEx transfer rate plots

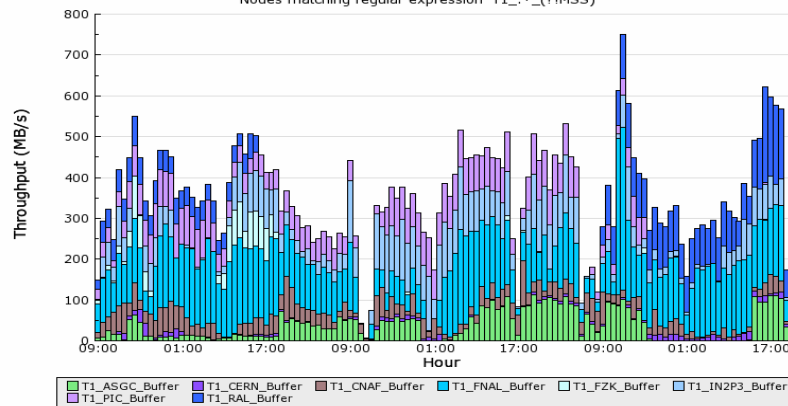
PhEDEx Prod Data Transfers By Destination

30 Days from 2006-10-02 to 2006-10-31 GMT
Nodes matching regular expression 'T1_.*(?:IMSS)'



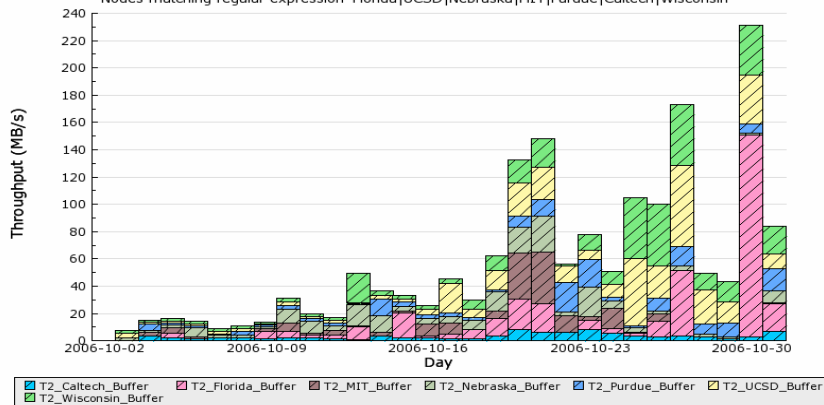
PhEDEx Prod Data Transfers By Destination

132 Hours from 2006-10-26 09:00 to 2006-10-31 20:00 GMT
Nodes matching regular expression 'T1_.*(?:IMSS)'



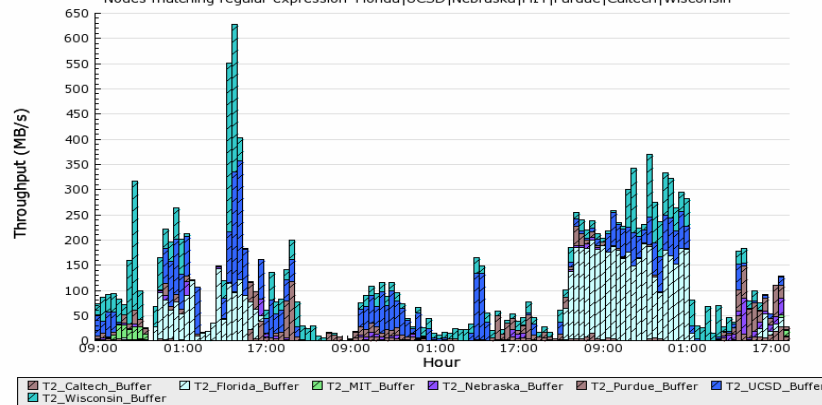
PhEDEx Prod Data Transfers By Destination

30 Days from 2006-10-02 to 2006-10-31 GMT
Nodes matching regular expression 'Florida|UCSD|Nebraska|MIT|Purdue|Caltech|Wisconsin'



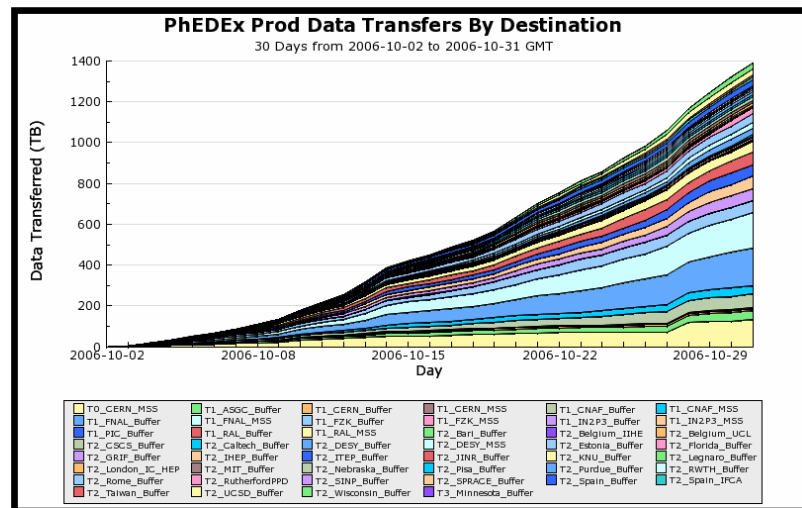
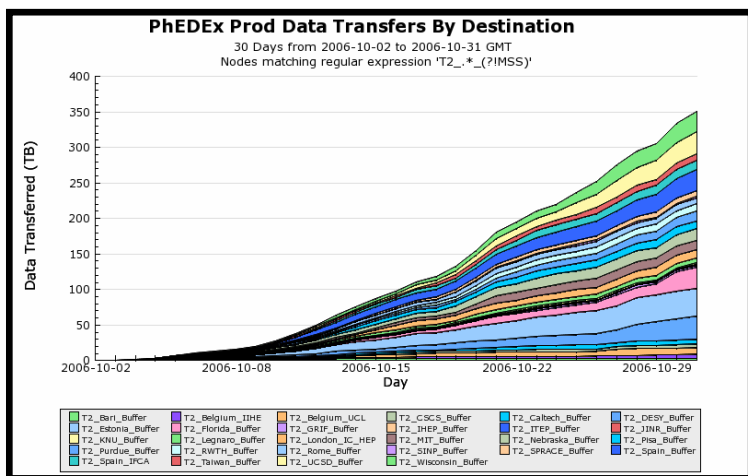
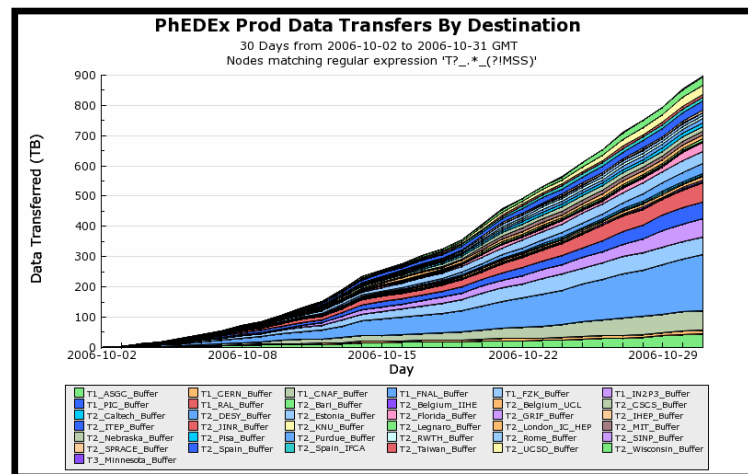
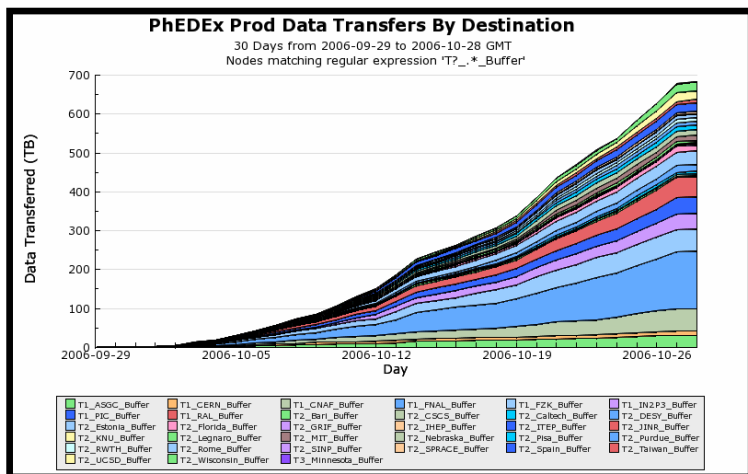
PhEDEx Prod Data Transfers By Destination

132 Hours from 2006-10-26 09:00 to 2006-10-31 20:00 GMT
Nodes matching regular expression 'Florida|UCSD|Nebraska|MIT|Purdue|Caltech|Wisconsin'





PhEDEx Data Volume





Goals and Metrics for CSA06

- Resources used for CSA06
 - About 1200 CPU at the Tier0
 - About 2000 to 2500 CPUs across all Tier1s
 - About 2000 to 2500 CPUs across all Tier2s
 - From 70 - 200TB disk + tape at participating Tier1s
 - From 5 - 25TB disk at participating Tier2s



MC Samples Created for CSA06

1. Minimum bias (40M)
2. T-Tbar (6M)
3. $Z \rightarrow \mu\mu$ (2M)
4. $W \rightarrow e\nu$ (4M)
5. Jet calibration soup (1M)
dijet + Z+jet, various pt-hat ranges
6. Electroweak soup (5M)
 $W \rightarrow l\nu$ + Drell-Yan ($m > 15$ GeV) + WW + H \rightarrow WW
7. Soft Muon Soup (2M)
Inclusive muons in minbias + J/Psi production
8. Exotics Soup (1M)
LM1 SUSY, Z' (700 GeV), and excited quark (2000 GeV)
9. HLT soup (5M)
W (leptons) + Drell-Yan (leptons) + t-tbar (all modes) + dijets
To be split into individual datasets for input to Tier0

**For calibration
exercises**

Total: ~66M events
Simulated up to
detector digitization
No pile-up



Data Placement

Centre	Minbias	T-Tbar	Z->mumu	W->enu	Jet Soup	EWK Soup	Soft Muon	Exotics Soup	HLT Streams
ASGC	10% (8TB)					10TB			
CNAF	15% (12TB)		4TB	8TB	2TB	10TB	4TB		
FNAL	35%(28TB)	12TB		8TB	2TB	10TB		2TB	10TB
GridKa	15%(12TB)	12TB				10TB		2TB	
IN2P3	15% (12TB)				2TB	10TB		2TB	10TB
PIC	5% (4TB)		4TB			10TB	4TB		
RAL	5% (4TB)	12TB				10TB		2TB	10TB



HLT Filters

<u>Name</u>	<u>Mnemonic</u>	<u>threshold</u>	<u>Bit position in 0 8 4</u>	<u>Bit Position 0 8 3</u>	<u>HLT Soup Efficiency</u>
Single Gamma	p1g	80	0	6	0.9%
Double Gamma	p2g	30,20	1	2	3.0%
Single electron	p1e	26	2	4	32%
Double electron	p2e	12,12	3	0	3.3%
Single Muon	p1m	19	4	9	35%
Double Muon	p2m	7,7	5	7	3.2%
Single Tau	p1t	100	6	12	0
Double Tau	p2t	60,60	7	11	0
Single Jet	p1j	400	8	8	2.7%
DiJet ²	p2j	350	9	10	2.2%
TriJet ²	p3j	195	10	1	0.5%
Quad Jet	p4j	80	11	5	0.6%