

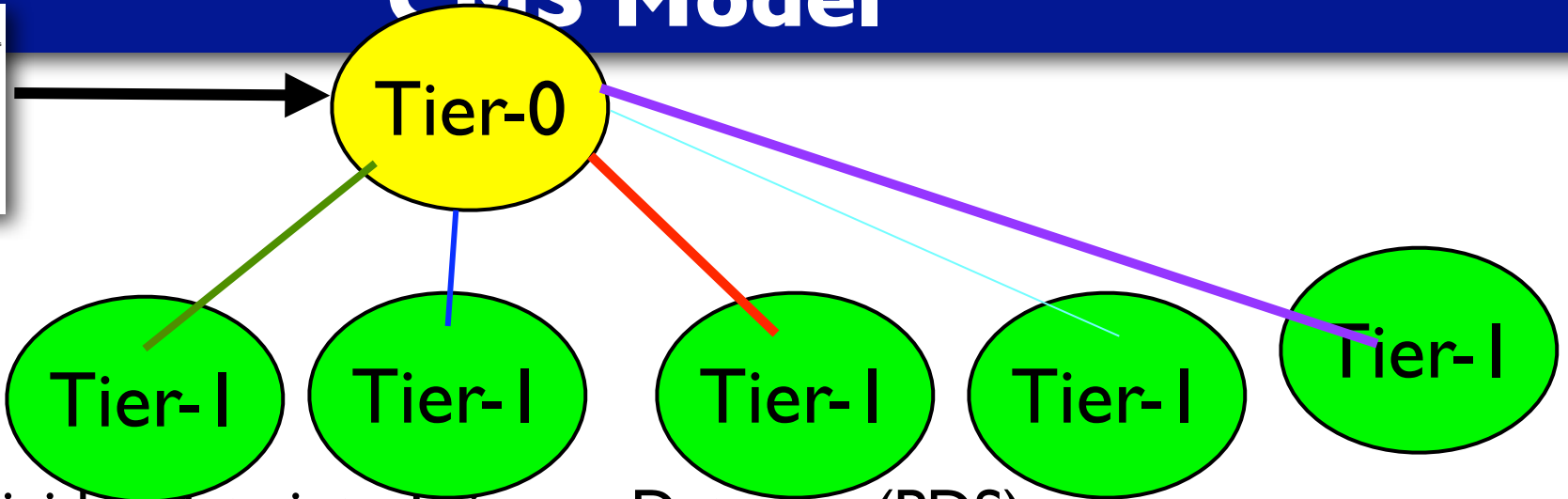
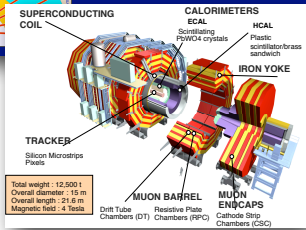
CMS Computing Experience

6 April 2010
DOSAR
Ian Fisk





CMS Model

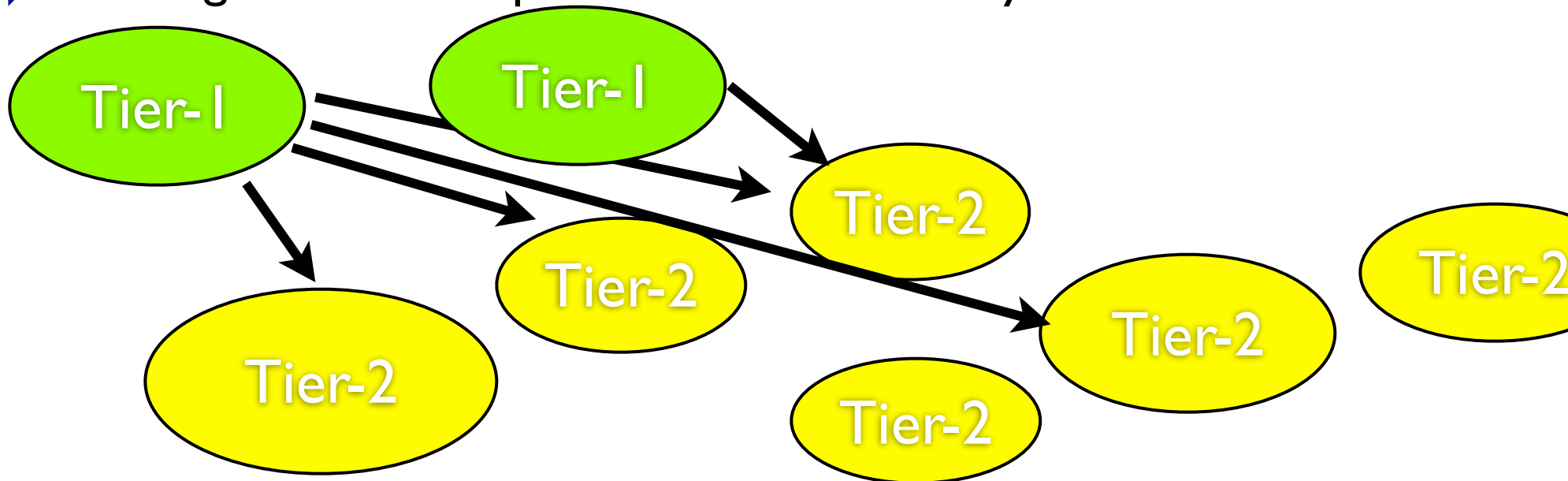


- ▶ CMS Divides data into Primary Datasets (PDS)
 - ▶ Currently we have 2
 - ▶ Growing to ~10 in 2010
 - ▶ Hopefully ~15 in 2011
- ▶ Datasets are based on trigger bits and will not all be the same size
 - ▶ Argued that the largest for sustainable long term operations is 50Hz (Out of 300Hz) and the smallest is 3-5Hz
- ▶ Datasets will also be cut in time
 - ▶ As data becomes incompatible or logically disconnected will stop and change datasets



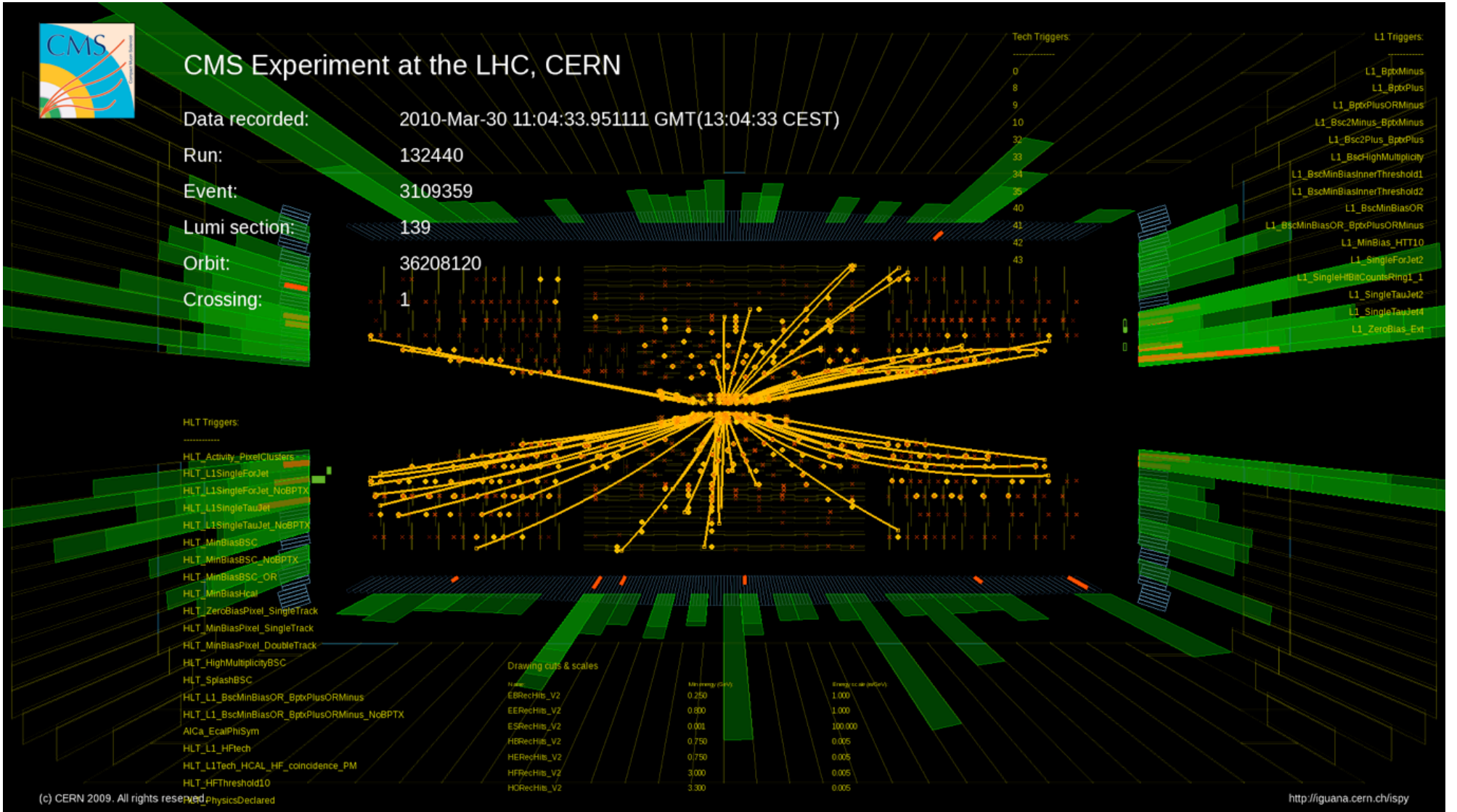
Exports from CMS

- ▶ Tier-I Centers are entrusted with CMS custodial data for archiving, processing, and data serving
- ▶ The Tier-2 sites in CMS requesting data samples can come from anywhere
- ▶ The rate to Tier-2s is driven by user bursts and is expected in the CTDR to reach 50MB/s for slower links to bursts to 500MB/s for faster links
- ▶ Serving of AOD samples can come from any Tier-I





7TeV Events





7 TeV Events



CMS Experiment at the LHC, CERN

Data recorded: 2010-Mar-30 11:04:37.067645 GMT(13:04:37 CEST)
 Run: 132440
 Event: 3111007
 Lumi section: 139
 Orbit: 36243167
 Crossing: 1

HLT Triggers:

-
- HLT_Activity_PixelClusters
- HLT_Activity_EcalREM
- HLT_L1Jet6J
- HLT_L1SingleForJet
- HLT_L1SingleForJet_NoBPTX
- HLT_L1SingleTauJet
- HLT_L1SingleTauJet_NoBPTX
- HLT_MinBiasBSC
- HLT_MinBiasBSC_NoBPTX
- HLT_MinBiasBSC_OR
- HLT_MinBiasHcal
- HLT_ZeroBiasPixel_SingleTrack
- HLT_MinBiasPixel_SingleTrack
- HLT_MinBiasPixel_DoubleTrack
- HLT_SplashBSC
- HLT_L1_BscMinBiasOR_BptxPlusORMinus
- HLT_L1_BscMinBiasOR_BptxPlusORMinus_NoBPTX
- HLT_L1_HFtech
- HLT_L1Tech_HCAL_HF_coincidence_PM
- HLT_HFThreshold10

Drawing cuts & scales

Name	Min energy (GeV)	Energy scale (GeV)
EBRecHit_V2	0.250	1.000
EERecHit_V2	0.800	1.000
ESRecHit_V2	0.001	100.000
HRRecHit_V2	0.750	0.005
HERecHit_V2	0.750	0.005
HFRRecHit_V2	3.000	0.005
HORecHit_V2	3.300	0.005

(c) CERN 2009. All rights reserved. Physics Declared

<http://figuana.cern.ch/ispv>



CMS Computing Model

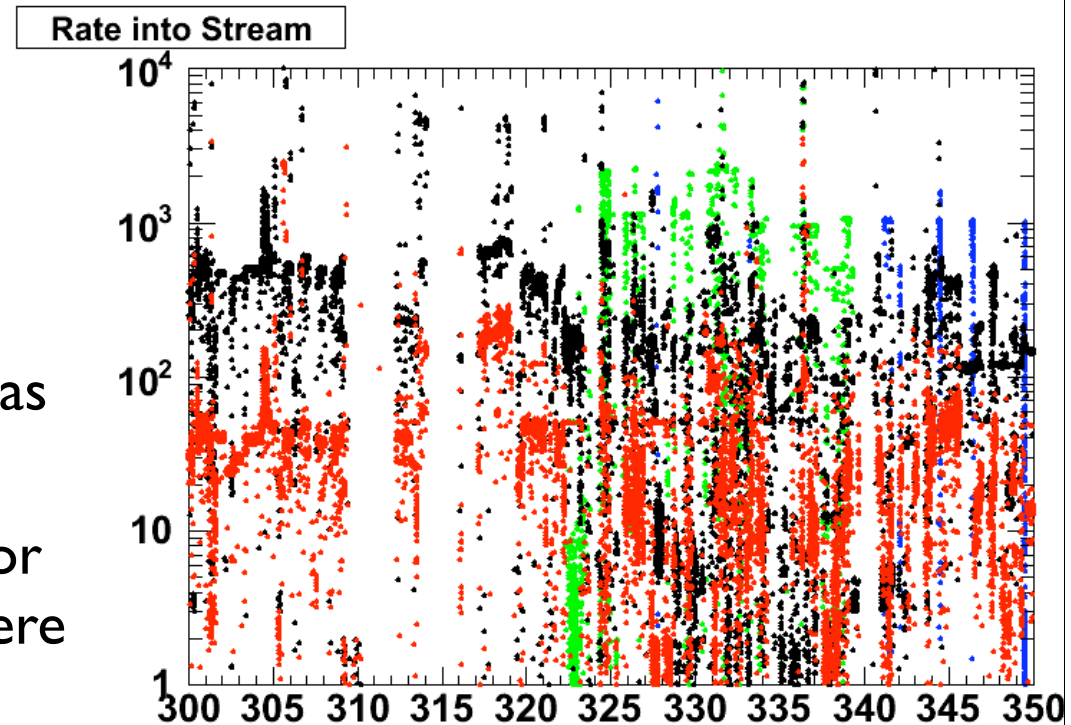
- ▶ Up to now the data period CMS performed workflows and activities that were predicted in the Computing Model
- ▶ Computing Tiers performed the specified workflows
 - ▶ Tier-0 promptly reconstructed, the Tier-1s reprocessed and served data, and the Tier-2 centers were used for simulated event production and analysis
 - ▶ An additional online stream was added
- ▶ The predicted workflows were executed much more frequently.
 - ▶ Reprocessing and analysis were exercised frequently
 - ▶ Data was subscribed to many T1s and more T2s
 - ▶ Replication and processing went well
 - ▶ Event complexity and fraction of “interesting” physics events was much lower than expected in the planning



Planning/Observation

- ▶ The Computing Model planning defined 2 periods
 - ▶ Period 1 is Oct-November 2009 to April of 2010
 - ▶ Period 2 was the remainder of 2010
- ▶ In Period 1
 - ▶ 100 days at 20% livetime (20 days)
 - ▶ 1.5MB/event RAW and 0.5RECO
 - ▶ Total Number of Events 726M
 - ▶ Total Volume of Data
 - ▶ ~1PB RAW
 - ▶ 359TB Reco
 - ▶ A few 10pb-1
 - ▶ Rate of Data from P5 450MB/s
- ▶ 2009 to Present for the Minimum Bias Sample
 - ▶ There are nearly 16k lumi sections on the RAW Minbias PDS
 - ▶ 17 days (90M events)
 - ▶ 2400 Files
 - ▶ 7.8TB
 - ▶ 10 inverse micro-barn collected
 - ▶ Stable beams with all detectors on and timed-in is smaller
 - ▶ 22 hours
 - ▶ 6.8M events (Around a 1TB)
 - ▶ 500k real collision events

- ▶ Data in CMS is sent from the Online as streams
- ▶ **Express** is expected to be about 40Hz. Generally stayed within 40-60Hz with occasion spikes to 3kHz
- ▶ **Stream A** is the source of the Primary Datasets (PDS). In the planning was expected at 300Hz, was 200Hz with spikes to 1kHz.
- ▶ Expectation in the planning was for 10 PDS. In the first run there were only 2 populated.
- ▶ **Stream B** was proposed before the run as insurance. It's a very high rate stream of ZeroBias Data. Averages 1kHz. **Stream B** was also buffered





Data volume: Streams

Stream	#Events	Size [GB]
HLTMON	19,454,692	3,935.81
Express	80,478,349	12,335.44
B_Buffered	130,167,201	25,478.95
A	731,269,373	98,467.30
B	278,019,843	20,111.04
RPCMON	145,150,042	540.61
FEDErrors	457	0.17
Calibration	209,228,981	24,186.10
ALCAP0	40,154,649	401.21
ALCAPHISYM	253,569,603	2,488.50
OnlineErrors	89,297	24.87
	1,887,582,487	187,970.01

- ▶ Planning for period I called for about 725M events
- ▶ 770M Simulated
- ▶ The corresponding data volume per month was IPB over 6 months
- ▶ Event size and complexity of processing much lower than planned
- ▶ The fraction of interesting to taken events much much lower



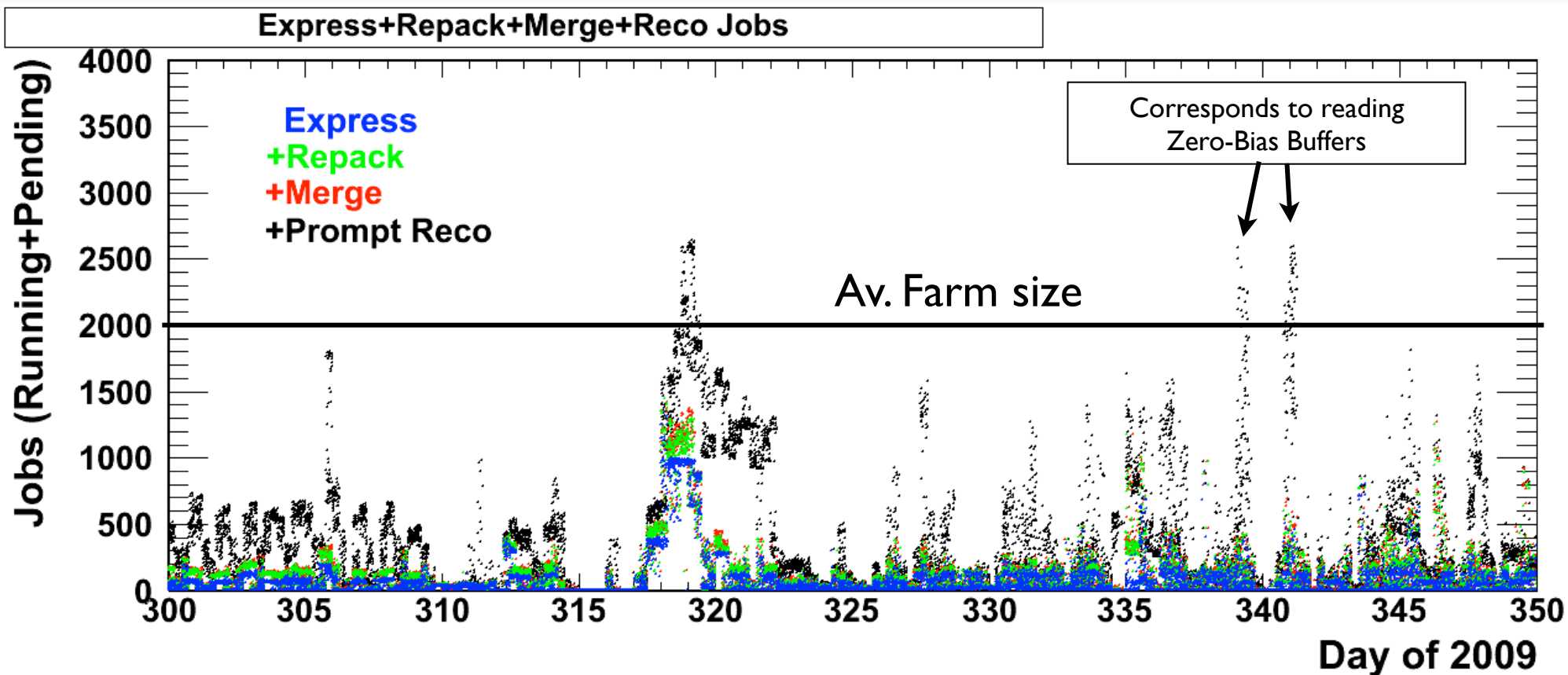
Data volume: PDs

PD	RAW		Prompt Reco	
	#Events	Size [GB]	#Events	Size [GB]
MinimumBias	90,052,125	7,822.32	89,791,258	6,964.36
RandomTriggersO	46,969	0.04		
ρ				
RandomTriggers	47,444,572	5,441.41		
ZeroBias	78,065,537	6,445.00	78,038,521	4,671.24
ZeroBiasB	404,057,754	37,999.02	20,124,018	1,389.00
LogMonitor	86,462	19.61		
TestEnables	209,228,219	19,352.02		
PhysicsMuonBkg	91,890,670	9,411.89	67,388,127	4,536.38
BeamHalo	123,852	17.40	80,971	31.21
Test	763,109	147.33		
AlCaPhiSymEcal	253,569,603	1,896.82		
MinimumBiasNoCa	0	0.02		
l				
AlCaP0	40,154,649	328.85		
HcalHPDNoise	1,674,393	257.93		
RPCMonitor	145,150,042	219.71		
FEDMonitor	3,293	0.39		
Calo	117,967,688	14,634.55	81,452,296	4,428.22
Cosmics	407,437,569	42,558.37	363,689,277	31,803.55
HcalNZS	6,916,109	1,203.36	61,619	8.77
ZeroBiasBnotT0	4,129,290	473.47		
	1,898,761,905	148,229.52	700,626,087.00	53,832.73

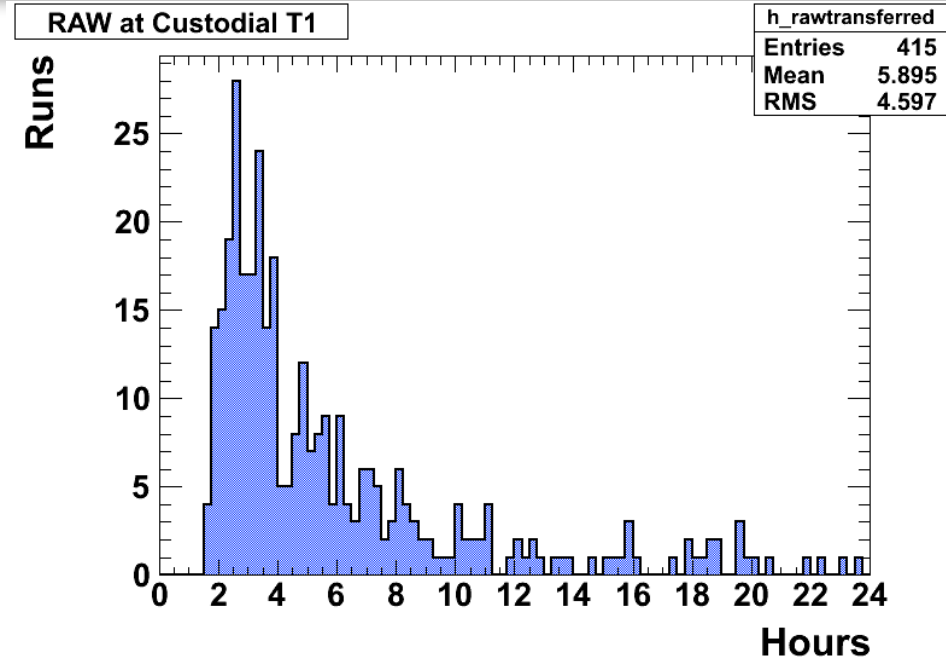
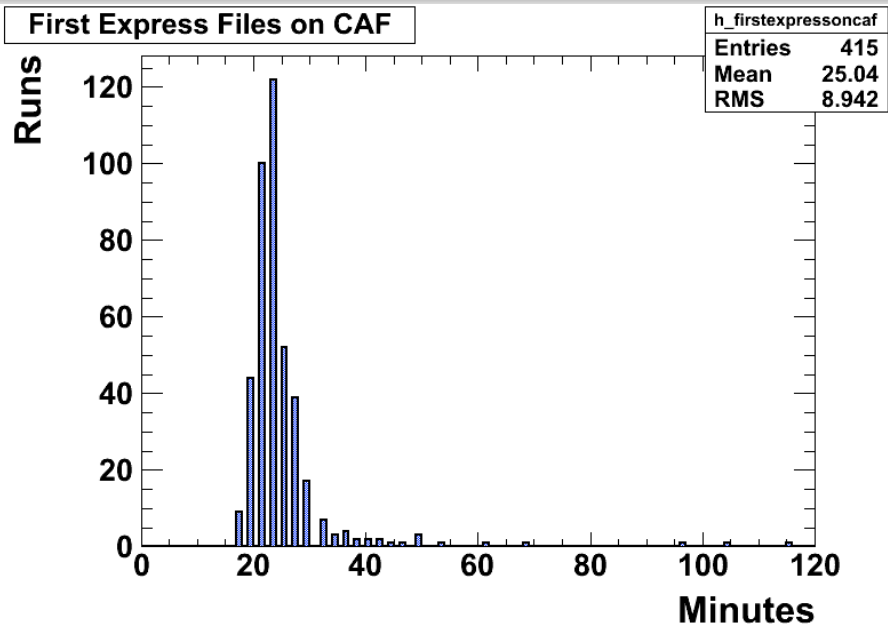
Average Event Size in Minimum Bias and Zero Bias Stream: 100k

► Agreement between real and Simulated event size and processing time for 900GeV is excellent

T0 queue utilization



- ▶ Each of the points are cumulative: black corresponds to all jobs
- ▶ Includes pending jobs (only jobs released from queue recorded)



- ▶ Time to first express files on the CAF relative to receiving first streamers of run at T0:

- ▶ The design spec for this time is one hour.
- ▶ Average more 25 min. With very small tails

- ▶ Latency for RAW data successfully transferred to the Tier-I for custodial storage:

- ▶ Average is within the model expectations
- ▶ The tails are understood and most are caused by specific technical issues



Data Distribution

BeamCommissioning09

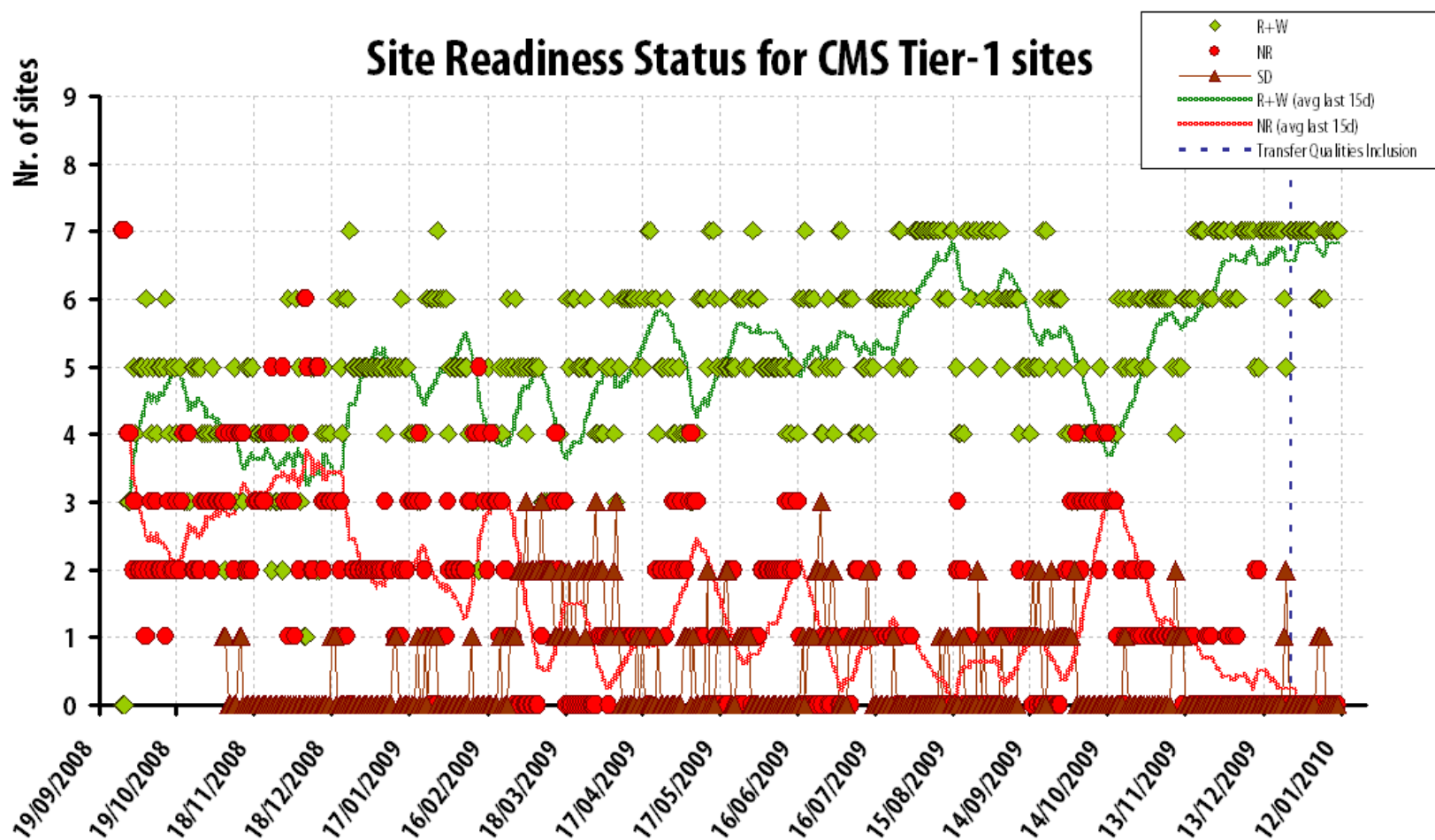
Primary Dataset	Custodial Site	Other Sites	RecoSequence	Comment
AlCaP0	RAL	FNAL	dedicated ALCARECO	-
AlCaPhiSymEcal	RAL	FNAL	dedicated ALCARECO	-
Cosmics	FNAL	RAL	cosmics	-
FEDMonitor	PIC	FNAL	None	-
HcalNZS	FNAL	-	collision+hcalnzs	-
HcalHPDNoise	FNAL	-	None	-
LogMonitor	FNAL	-	None	-
MinimumBias	IN2P3	FNAL, PIC, CNAF	collision	-
PhysicsMuonBkg	FNAL	IN2P3	collision	new
RandomTriggers	PIC	FNAL	None	-
RPCMonitor	PIC	FNAL,CNAF	None	-
TestEnables	FNAL	-	None	-
ZeroBias	CNAF	FNAL	collision	new - 100 Hz in A stream
ZeroBiasB	KIT	FNAL	collision or None depending on resources	new - 1 kHz in B stream
ZeroBiasBnotT0	RAL	-	None	new

► MinimumBias was replicated to 4 TI sites in total



Tier-1 Readiness

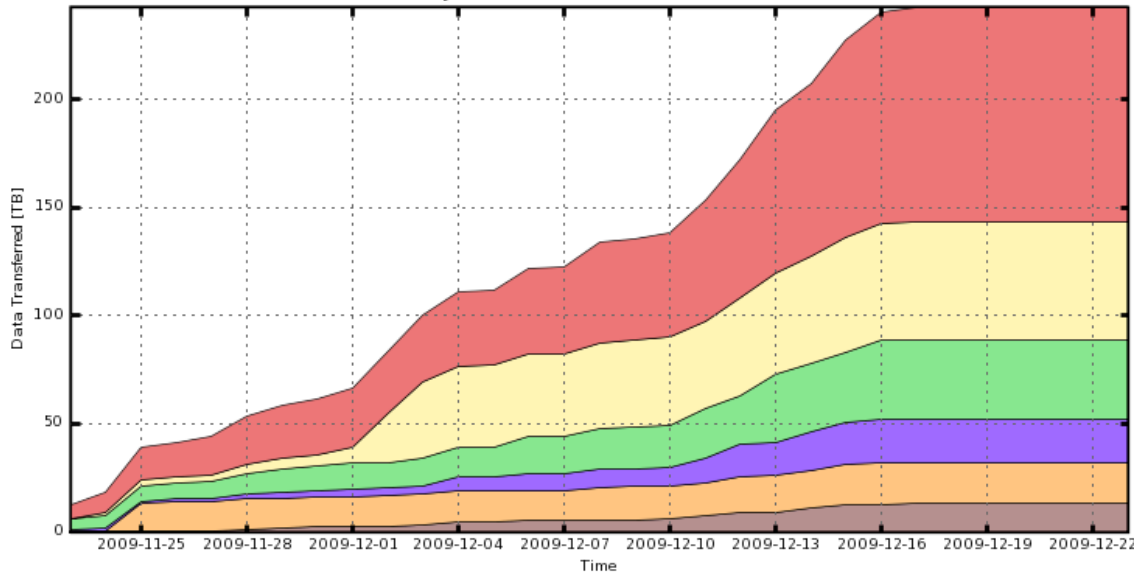
- ▶ Site readiness of the Tier-1s has improved, but CMS ran with only 6 sites receiving custodial data during the 2009 collision data
- ▶ Goal is to send data to 2 sites plus an archival copy at CERN for as long as the resources permit in 2010





Transfers: T0 → T1

CMS PhEDEx - Cumulative Transfer Volume
30 Days from 2009-11-23 to 2009-12-23



■ T1_US_FNAL_Buffer
 ■ T1_FR_CCIN2P3_Buffer
 ■ T1_DE_KIT_Buffer
 ■ T1_IT_CNAF_Buffer
 ■ T1_UK_RAL_Buffer
■ T1_ES_PIC_Buffer

Total: 242.43 TB, Average Rate: 0.00 TB/s

Site	Total Transfer Volume [TB]	Percentage	Expected
T1_DE_KIT	36.68	15	12
T1_ES_PIC	13.03	5	7
T1_FR_CCIN2P3	54.33	22	12
T1_IT_CNAF	20.23	8	8
T1_UK_RAL	19.08	8	10
T1_US_FNAL	99.07	41	40
	242.42	99	89

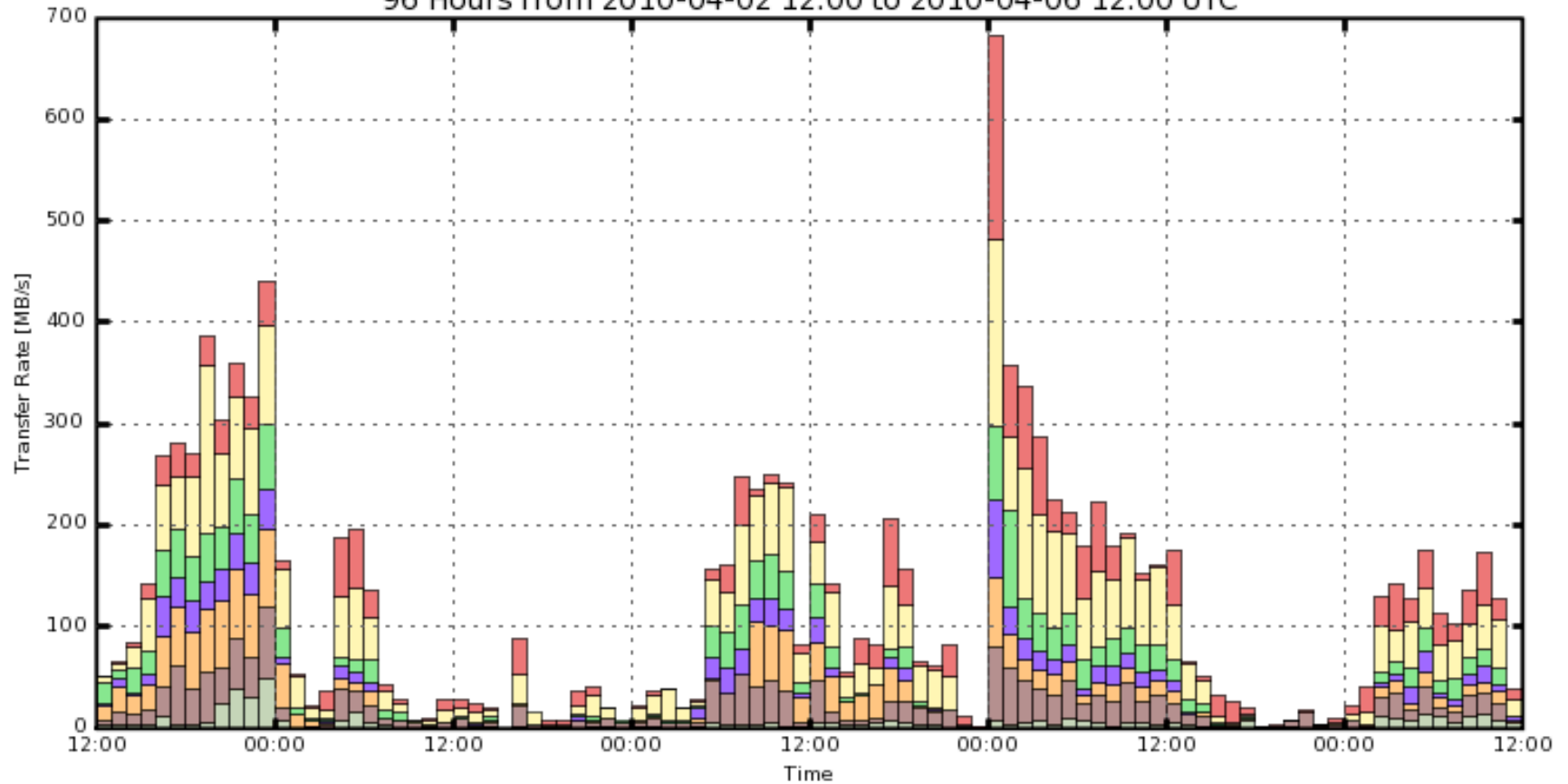
Doesn't include ASGC

- ▶ This only shows data originating at T0
- ▶ 6 T1 sites received data



CERN to Tier-1 Transfers

CMS PhEDEx - Transfer Rate
96 Hours from 2010-04-02 12:00 to 2010-04-06 12:00 UTC



T1_IT_CNAF_Buffer T1_US_FNAL_Buffer T1_DE_KIT_Buffer T1_ES_PIC_Buffer T1_UK_RAL_Buffer
T1_FR_CCIN2P3_Buffer T1_TW_ASGC_Buffer

Maximum: 682.47 MB/s, Minimum: 0.00 MB/s, Average: 122.68 MB/s, Current: 39.17 MB/s



TI processing

- ▶ 5 rereco passes so far (not counting the first small one)

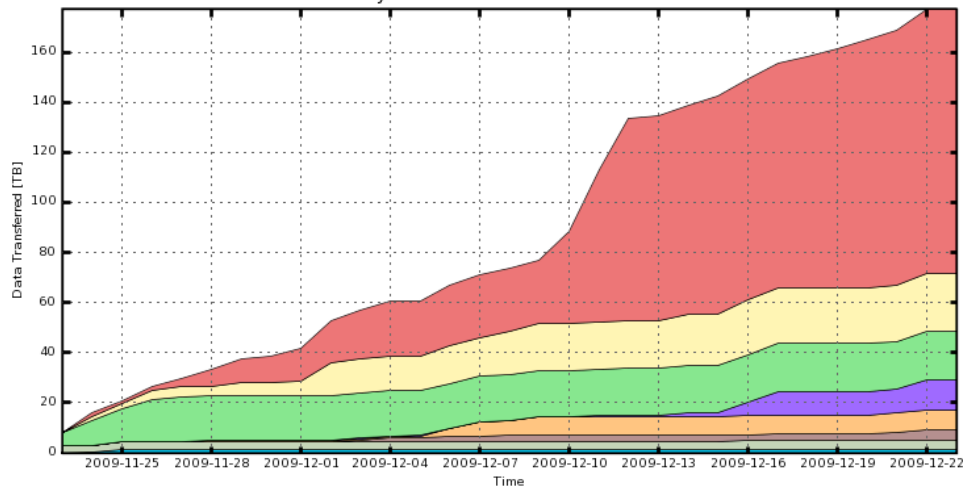
Rereco Pass	#Runs	#Events	#Input-Files = #Jobs	Rough Expectations
First rereco (GR09_P_V7)	11	998,545	76	
Dec09	14	3,499,729	144	
Dec14	37	15,885,562	701	
Dec19	52	19,681,382	925	250M Events

- ▶ total number of produced events: 150M
- ▶ total output size: 34TB (includes all data tiers)
- ▶ Latency: ~1-2 days (Planning expectations 1-2 weeks)
- ▶ Main time consumption:
 - ▶ Long running jobs (many events in input file while splitting by file to keep lumi sections intact)
 - ▶ Debugging and bookkeeping
- ▶ CPU efficiency: ~80-90% for reprocessing jobs
- ▶ Still some errors in monitoring and memory applications



Transfers: T1->T1

CMS PhEDEX - Cumulative Transfer Volume
30 Days from 2009-11-23 to 2009-12-23

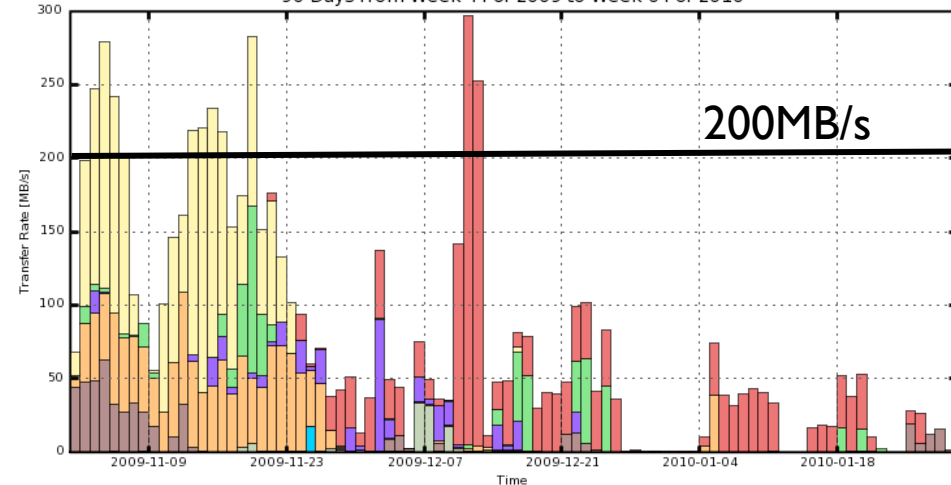


Total: 176.92 TB, Average Rate: 0.00 TB/s

(*) includes 2.93 TB transfers from TI_DE_FZK to TI_DE_KIT, 7.51 TB to repair samples at ASGC, 23.31 TB going to TI_CH_CERN

Destination Site	Total Transfer Volume [TB]
T1_DE_KIT	0.39
T1_ES_PIC	1.51
T1_FR_CCIN2P3	105.15
T1_IT_CNAF	4.55
T1_UK_RAL	19.27
T1_US_FNAL	12.29
Total	143.16

CMS PhEDEX - Transfer Rate
90 Days from Week 44 of 2009 to Week 04 of 2010



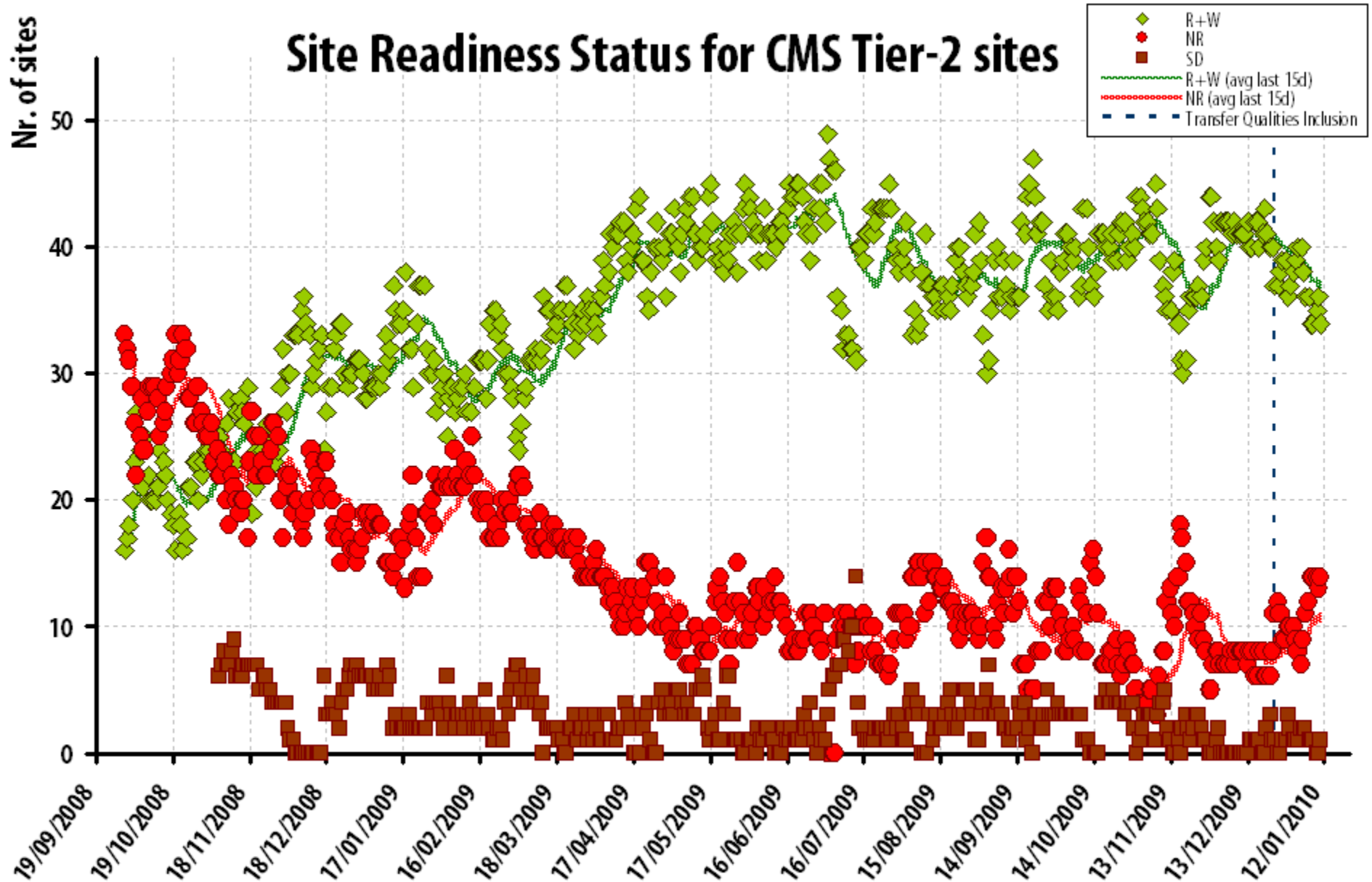
Maximum: 297.19 MB/s, Minimum: 0.00 MB/s, Average: 76.30 MB/s, Current: 1.03 MB/s

► Data transfer between T1 sites dominated by repopulation of IN2P3



Tier-2 Readiness

► Tier-2 Readiness has plateaued

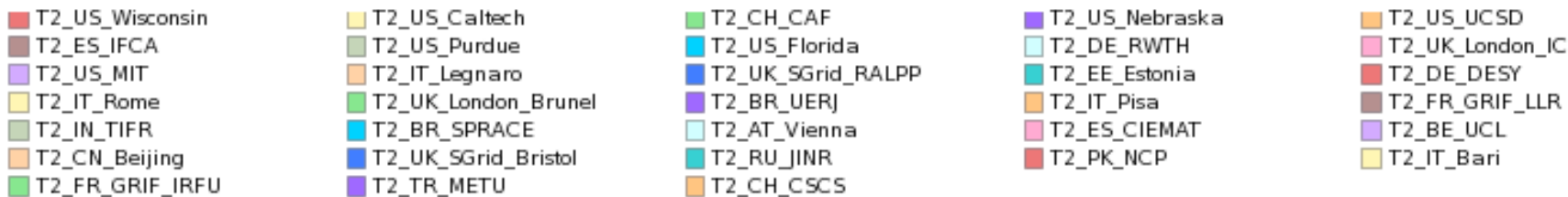
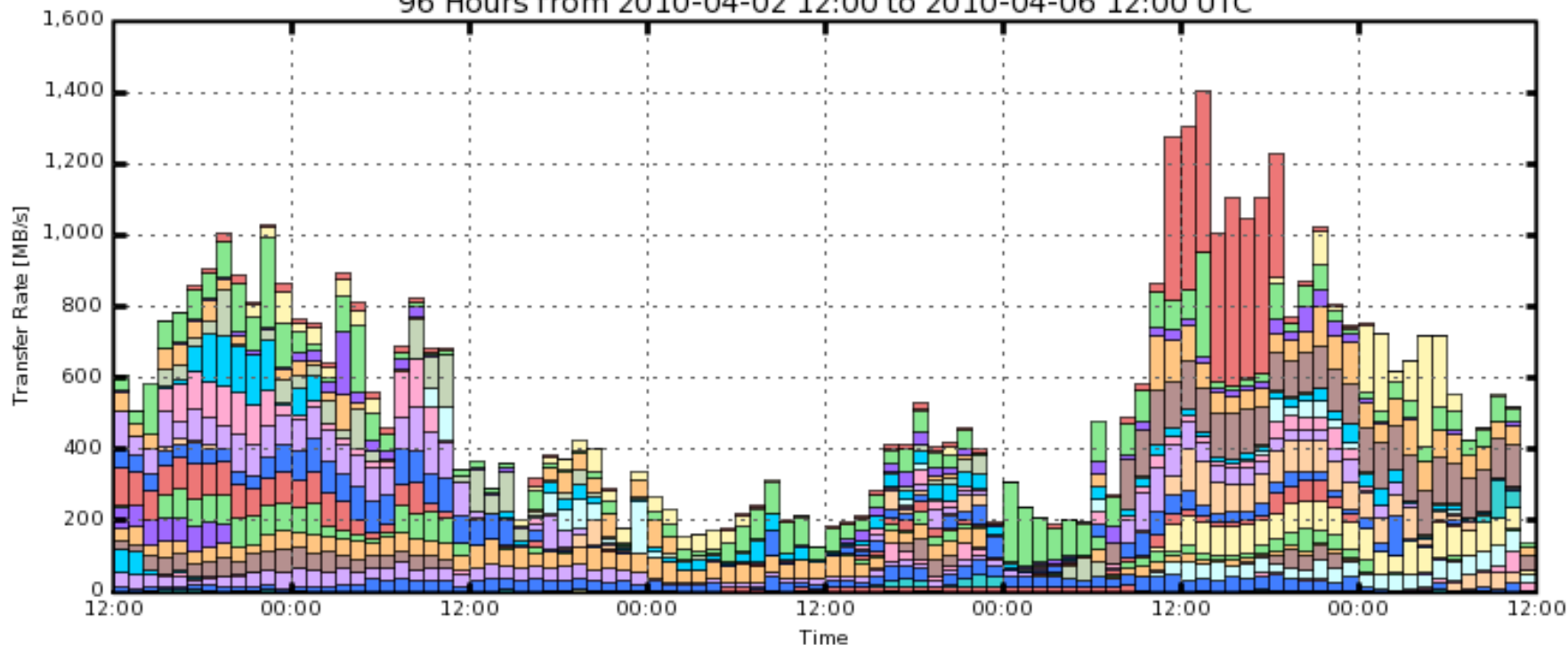




Transfers to Tier-2s

CMS PhEDEx - Transfer Rate

96 Hours from 2010-04-02 12:00 to 2010-04-06 12:00 UTC

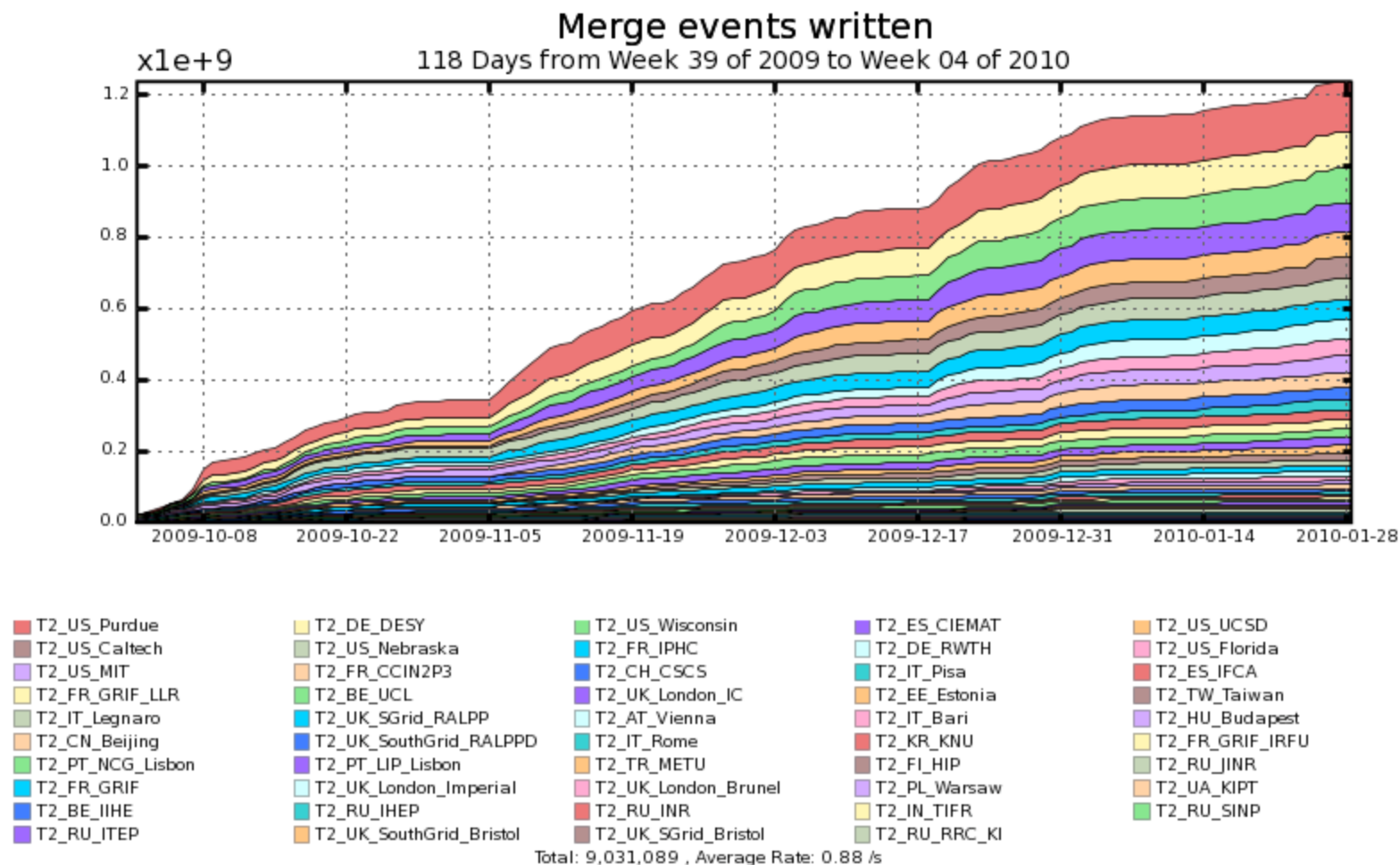


Maximum: 1,401 MB/s, Minimum: 126.34 MB/s, Average: 549.59 MB/s, Current: 137.36 MB/s



MC production

- ▶ Planning period I started in Oct
- ▶ 1.2B Events = ~ 400M individual simulation events roughly scales where we expected to be
- ▶ 3 months through 6 month period we have half of 780M

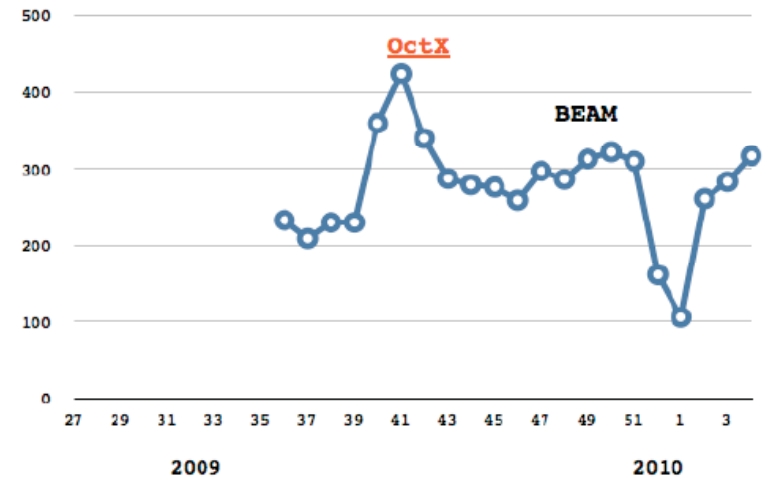




Analysis

- ▶ Since last review CMS has formed an Analysis Operations Team
 - ▶ Provide technical support for analysis infrastructure
 - ▶ Subscribe samples to centrally controlled space at Tier-2s
 - ▶ Analysis Ops has access to 50TB of space at ~50 sites
 - ▶ Currently ~1PB of space is utilized

Number of Analysis Users at Tier-2 Sites Each We

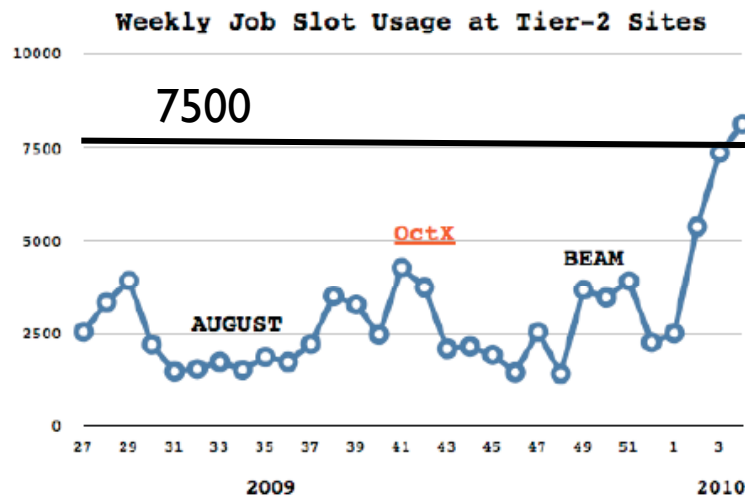


- ▶ Total number of people submitting distributed analysis jobs in a given week ~300
- ▶ Bump after the October analysis exercise

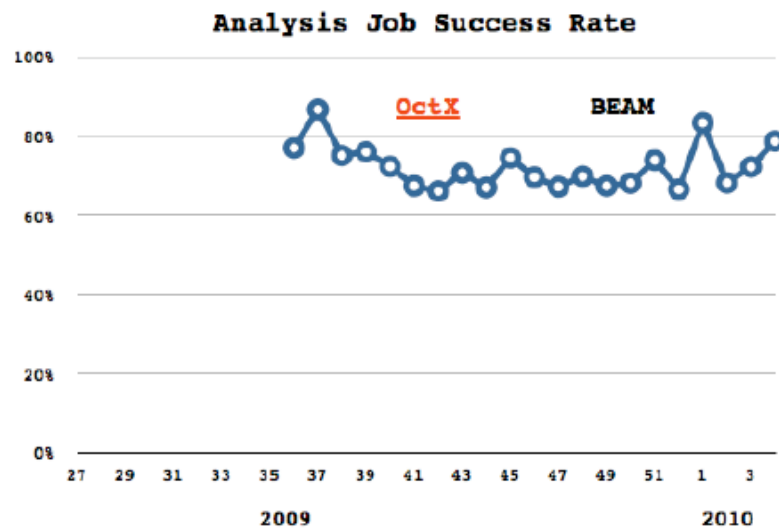
Type of Data	Total Size (TB)	All Replicas		AnalysisOps Requested	
		Size (TB)	Sites	Size (TB)	Sites
Collision Data	446.9	688.8	29	218.0	12
900 GeV MC	44.8	159.8	32	70.7	14
2.36 TeV MC	31.1	72.6	24	48.7	10
All Data Sets	9.3 PB		50	1.04 PB	38



Analysis



- ▶ Roughly 11k jobs slots are available for analysis
- ▶ Reaching 75% utilization toward the beginning of the year
- ▶ In any given week 47+/-2 Tier-2 sites have analysis jobs



- ▶ Success rate remains a persistent issue
- ▶ Improvement over last year where we had ~65%
- ▶ Half of errors are related to remote stage-out of produced files



Utilization

- ▶ Activities like MC and Analysis that are driven by external factors are making reasonably high use of the available resources
 - ▶ Analysis is currently running at 75% level
 - ▶ MC roughly on planning
- ▶ Activities like Re-reconstruction and skimming that are driven by available data are not fully utilized
 - ▶ Data Volume is much lower than planned for
 - ▶ Transfers lower on average. Good peaks. Partially compensated by over subscription
 - ▶ Tier-I utilization for activities like Cosmic reprocessing was high. On average lower than planned.



Outlook

- ▶ CMS is looking at approximately 1M times more integrated luminosity by the start of the summer
 - ▶ 10 inverse micro-barn to 10 pb⁻¹
- ▶ While many elements of the computing model accurately reflect the activities and the experience,
 - ▶ We have little experience with a resource constrained system
 - ▶ We don't have experience with large quantities of very interesting physics events
 - ▶ We are hoping for a huge increase in data volume.