



# CMS Service Challenge Plans

Ian Fisk  
June 21, 2006



# Service Challenge Plans

From the standpoint of CMS computing, Service Challenge 4 is a dress rehearsal for CSA06

- ➔ A number of the elements in SC4 are lower scale exercises we expect to perform in CSA06
- ➔ There primary technical differences are
  - We add Tier-0 reconstruction (outside the scope of today)
  - We anticipate include user analysis jobs into Tier-1 and Tier-2 processing
    - Currently SC4 is entirely exercised with robots
    - Job submission rates for CSA06 are 50k jobs per day (double SC4)
  - Calibration and Alignment infrastructure plays a larger role
- ➔ SC4 is an integration activity we will continue running for as long as it continues to be productive
  - CSA06 is a data challenge with metrics that need to be met in a specified time



# Operational Goals

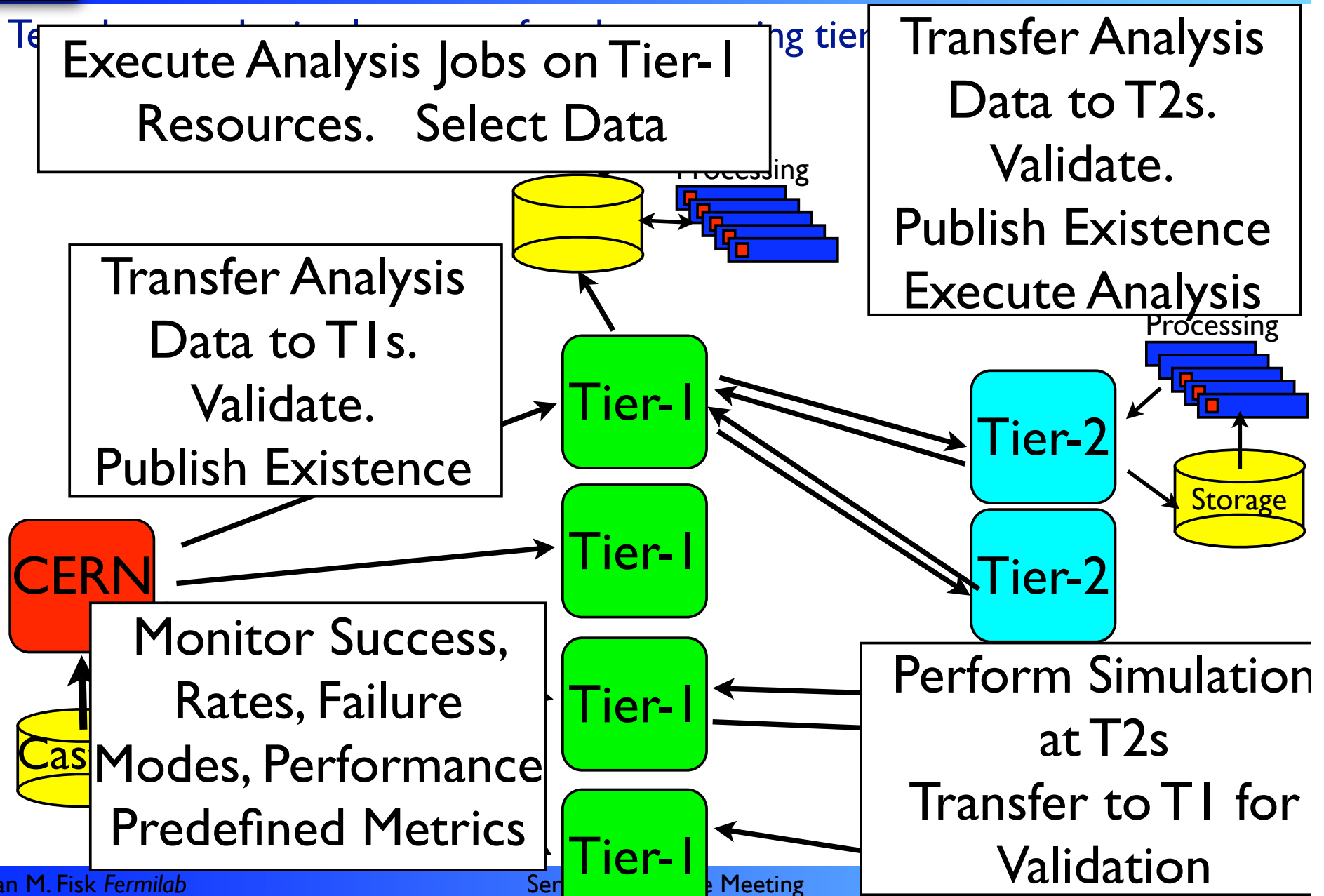
CMS needs to be at production scale services in 2008

- ➔ Assuming we cannot easily more than double the scale each year, we should be able to demonstrate 25% of the expected 2008 scale in this year and be able to reach 50% scale early in 2007

Service	2008 Goal	2006 Goal	%
Network Transfers between T0-T1	600MB/s	150MB/s	25%
Network Transfers between T0-T1	50-500 MB/s	10-100 MB/s	20%
Job Submission to Tier-1s	50k jobs/d	12k jobs/d	25%
Job Submissions to Tier-2s	150k jobs/d	40k jobs/d	25%
MC Simulation	1.5 $10^9$ events/year	25M per month	25%



# SC4 Workflows





# Original Schedule Processing

Original schedule was to operate the first two weeks of June

- ➔ 25k jobs per day (50% analysis and 50% production)
  - Operate Job Robot on test simulation samples for analysis
  - Operate Prod\_Agent for production job
- ➔ 90% success rate to complete jobs
- ➔ We had a series of validation steps for the sites to meet
  - 25 Tier-2s signed up to participate
    - Pass the site functional tests, allow the CMS software to be installed, configure PhEDEx, download a test sample into the trivial file catalog namespace, demonstrate access with an analysis application
    - Demonstrate success with the production agent job
  - 19 Tier-2 sites and 5 Tier-1 sites completed the steps
- ➔ We spent a lot of the first two weeks commissioning sites and did not start large scale operations until last week



# Original Schedule Transfers

## Scaling Tape Rates by pledge aiming for 150MB/s

- ➔ ASGC: 10MB/s to tape
- ➔ CNAF: 25MB/s to tape
- ➔ FNAL: 50MB/s to tape
- ➔ GridKa: 20MB/s to tape
- ➔ IN2P3: 25MB/s to tape
- ➔ PIC: 20MB/s to tape
- ➔ RAL: 10MB/s to tape

## Networking provisioning should be at least twice this

- ➔ These goals are sufficiently modest that no center should struggle to sustain them



# Tier-2 Transfers

## Network Estimates for Tier-2 vary widely

- ➔ The computing model defines the expected minimum in 2008 at 1 Gb/s
  - Naively taking 25% this would be 250Mb/s
- ➔ Given the number of Tier-2 centers already at 1 Gb/s to 10Gb/s and the difficulty using reasonable scale networking end-to-end it makes sense to try much larger scale tests at some Tier-2 centers
  - Try to sustain ingest rate to Tier-1 centers from all Tier-2s
  - Drive Tier-1 to Tier-2 rates at 10MB/s to 100MB/s



# Estimated Performance of Mass Storage

The expected mass storage performance

- ➔ At a nominal Tier-1 is 800MB/s to the worker nodes
- ➔ At a nominal Tier-2 is 200MB/s

We have been aiming for 1MB/s per batch slot

- ➔ 300MB/s at a nominal Tier-1
- ➔ 100MB/s at a nominal Tier-2

These numbers are higher than the simple 25% scaling, but should be well within capacity for reasonable size Tier-2 centers.

These are being exercised and documented in CSA06

- ➔ For Castor and dCache the performance has be good
- ➔ DPM is being exercised in SC4

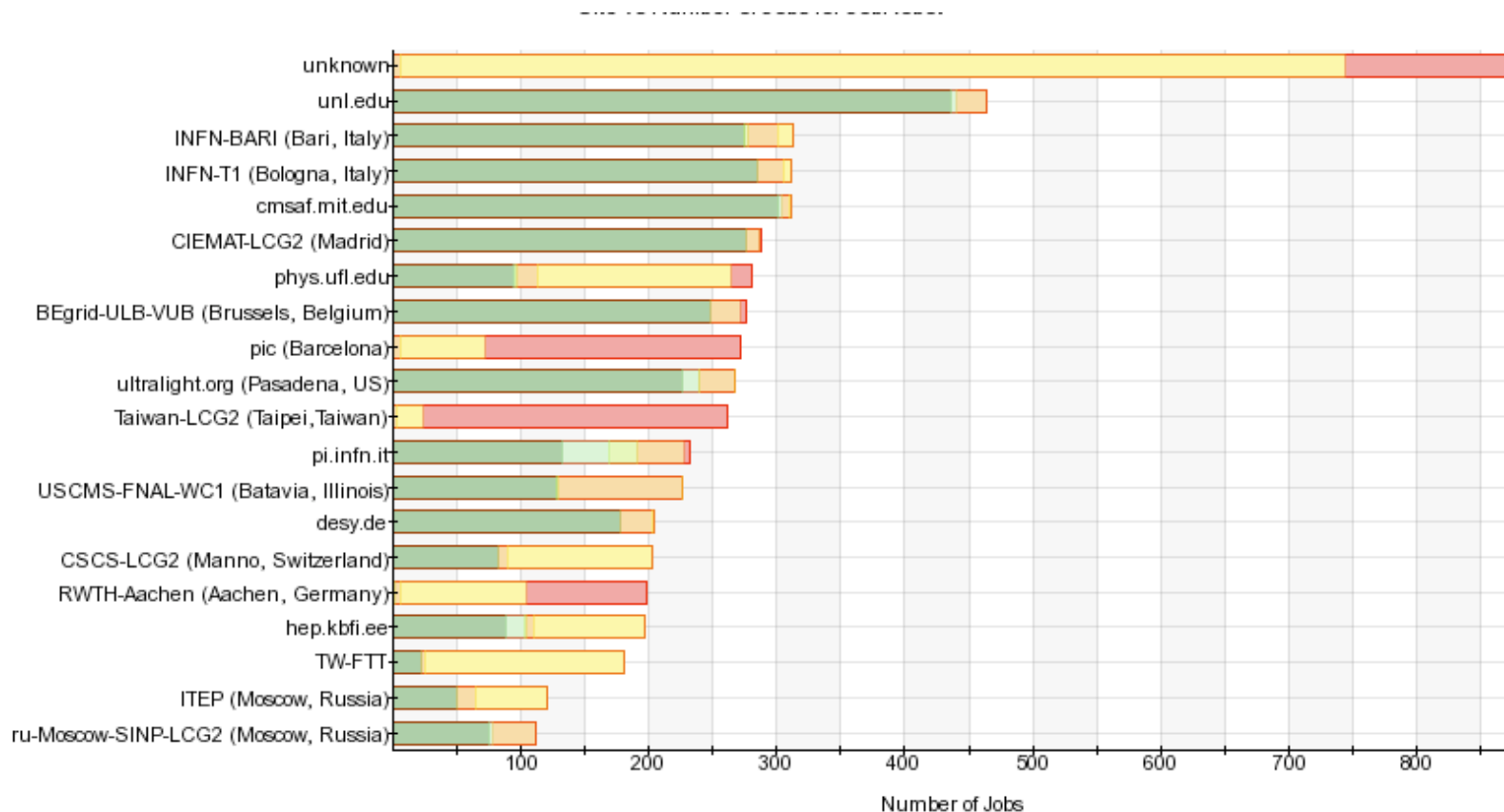




# Processing Status

Analysis processing is going OK, though we still have about a factor of 2 in scaling to meeting and some stability issues

- ➔ We have hit about 6k analysis submissions in a day
- ➔ We have hit 10k jobs a day in the production system on the OSG
- Lower on the LCG





# Transfer Status

## Transfers proceeding poorly

- ➔ Data rate out of CERN has struggled to reach 100MB/s and then with substantial numbers of errors and large structure
  - Periods with 10MB/s-20MB/s
  - Many errors
- ➔ Rate from Tier-1s to Tier-2s are not much better
  - Highest export rate is from FNAL that has not switched the Tier-2s to FTS channels

## We have a couple exceptions

- ➔ A number of Tier-2 have sustained 10MB/s from at least one Tier-1
- ➔ ~5 Tier-2s have sustained 50MB/s for a 24 hour period from one T1
- ➔ 1 Tier-2 has sustained 100MB/s for a 24 hour period



# Preparation Activities

## Increasing the processing scale and including users

- ➔ Need to hit 50k jobs per day by CSA06

## Continuing to exercise the transfers

- ➔ Get Tier-0 to Tier-1 transfers under control and begin to have more successful Tier-1 to Tier-2

## Incorporating the calibration information into the workflow

- ➔ Deploy the LCG-3D infrastructure with Frontier for SQUID caches

## Improving the functionality and reliability of sites

- ➔ All participating sites should be able to complete the CMS workflow and metrics

## Prepare CSA06 Production Samples

- ➔ 25M Events per month in July and August



# Increasing the Processing

## Two improvements from SC4

- ➔ To grow from 25k jobs per day to 50k jobs per day we need to switch submission infrastructure
  - 25k jobs per day is already a strain on the resource broker infrastructure we have
- ➔ Job robots are good load generators, but they do not make mistakes and they are patient (flat load over a 24 hour period)
  - Users generate unexpected usage patterns and loads

## The switch to the gLite RB with bulk submission is needed for CSA06

- ➔ Deployment came later than we expected in May
  - Still not fully commissioned in CMS for SC4
  - We anticipate the latter portion of SC4 and continuing throughput the summer to work with the gLite submission infrastructure on LCG
  - Continue to improve the submission rate on OSG



# Improving the Transfers

## Issues were recognized in the SC4 throughput tests

- ➔ It takes too long to get going and it takes too much effort to keep going
  - Services needed to be developed in real time
  - The response is to run the challenge more often until it becomes automatic
- ➔ In order to meet the transfer goals sites had to have a lot more transfers in action than the Computing models call for

## In SC4 CMS is observing equally important issues

- ➔ Rates achieved in the throughput tests have not be reproduced in the experiment tests.
- ➔ Experiment data management systems should drive transfers
  - At the end the raw FTS rate with a full time expert is not as important as the ability to move actual experiment data



# Proposed Extension

Starting in July WLCG would start working with ATLAS and the ATLAS data management system

➔ In the second half of July add CMS and PhEDEx

The total resource needs remain fairly manageable

<i>Centre</i>	<i>ATLAS</i>	<i>CMS*</i>	<i>Combined</i>	<i>SC4 proposal</i>	<i>Nominal</i>
ASGC	60.0	10	70	60	100
CNAF	59.0	25	85	60	200
PIC	48.6	30	80	50	100
IN2P3	90.2	15	115	100	200
GridKA	74.6	15	95	80	200
RAL	59.0	10	70	60	150
BNL	196.8	-	200	200	200
TRIUMF	47.6	-	50	50	50
SARA	87.6	-	90	90	150
NDGF	48.6	-	50	50	50
FNAL	-	50	50	-	200
Totals			~1000	800	1600



# Service Requirements

The ideas are still forming, but the proposal was to have an experiment coordinator (Someone to run the PhEDEx system) and a Service Coordinator (Someone to make sure transfer service is working)

- ➔ During the SC4 throughput phase this person worked extremely hard
- ➔ Proposal that the service coordinator is a large term need and people would take shifts

This is an interesting change of support philosophy

- ➔ The service coordinator model appears to acknowledge that currently the service requires some baby sitting at least one expert is there
  - The second philosophy appears to be much more realistic in the current level of service maturity

For CMS July is a reasonable time. The CPU resources will be tied up with simulation

- ➔ Network and storage should be available, need an experiment person



# Incorporating Calibration into Workflow

The Integration program has demonstrated the functionality of the Frontier system for distributing read-only database information since December

- ➔ CMS note explaining the very promising results

Frontier caches have been deployed at 10 sites for CMS testing

- ➔ Beginning in late June, all the participating CMS SC4 sites should get Frontier caches within the context of the LCG-3D

Site configuration system has been refined to include a discoverable local configuration for nearest database cache

Big item is to have real software applications that rely on the calibration information and prototype calibrations to populate the infrastructure

- ➔ Anticipated by the end of the summer
- ➔ We will continue to test with applications that simulate the access patterns until the real applications are available





# Improving Functionality and Reliability of Sites

Commissioning sites was a fair bit of work, but we are succeeding

- ➔ There are a number of new services and new sites

We have conducted and will continue open technical support sessions

- ➔ Opportunities for sites to call in and receive help with individual services

We have encouraged status reports

- ➔ A regular schedule of site reports
- ➔ Directed reports from Tier-I centers on SC4 progress

We will interact with sites through the LCG-MB

We are exercising the channels available to us, but there are still issues with site preparation and reliability

- ➔ The majority of sites are responsive, but there is a lot of work for this summer



# CSA06 Preparations

## The sample is 50M Events

- ➔ As a rough assumption let us assume 3 minutes per event on average
  - Based on the old GEANT4 production
- ➔  $50\text{M} \times 3 \text{ minutes} = 150\text{M} \text{ minutes}$
- ➔ July and August = 86k minutes
- ➔ If we were 100% efficient we would need 1750CPUs
  - Assuming something more reasonable like 75% we need ~2400
- ➔ Given these assumptions, most of the non-CERN resources needed for CSA06 itself will also be for pre-challenge production
- ➔ Requirements go down if CERN resources are used for simulation as well
  - Hoping to reserve CERN resources for Tier-0 preparations



# Rough Schedule

## Now until the end of June

- ➔ Continue to try to improve transfer efficiency
- ➔ Attempt to hit 25k jobs per day and increase the number and reliability of sites performing 90% efficiency for job completion

## July

- ➔ Demonstrate CMS analysis submitter in bulk mode with the gLite RB

## July and August

- ➔ 25M events per month with the production systems

## Second half of July participate in multi-experiment FTS Tier-0 to Tier-1 transfers

- ➔ Continue through August with transfers

## Improve Tier-1 to Tier-2 transfers and the reliability of the FTS channels.