# The LCG Service Challenges and GridPP

Jamie Shiers, CERN-IT-GD

31 January 2005

# Agenda

_LCG Project,  Grid Deployment Group, CERN_

- Reminder of the Goals and Timelines of the LCG Service Challenges

- Summary of LHC Experiments' Computing Models

- Outline of Service Challenges

- Review of SC1

- Status of SC2

- Plans for SC3 and beyond

# Why am I here?

**<u>As we will see in more detail later:</u>**

- Summer 2005: SC3 - include 2 Tier2s; progressively add more

- Summer / Fall 2006: SC4 complete

- SC4 – full computing model services
  - Tier-0, ALL Tier-1s, all major Tier-2s operational at full target data rates (~1.8 GB/sec at Tier-0)
  - acquisition - reconstruction - recording – distribution, *PLUS* ESD skimming, servicing Tier-2s

- How many Tier2s?
  - ATLAS: already identified 29
  - CMS: some 25
  - With overlap, assume some 50 T2s total(?)

- This means that in the 12 months from ~August 2005 we have to add 2 T2s per week
  - Cannot possibly be done using the same model as for T1s
    - SC meeting at a T1 as it begins to come online for service challenges
    - Typically 2 day (lunchtime – lunchtime meeting

# GDB / SC meetings / T1 visit Plan

- In addition to planned GDB meetings, Service Challenge Meetings, Network Meetings etc:

- Visits to all Tier1 sites (initially)
  - Goal is to meet as many of the players as possible
  - Not just GDB representatives! Equivalents of Vlado etc.

- Current Schedule:
  - Aim to complete many of European sites by Easter
  - "Round world" trip to BNL / FNAL / Triumf / ASCC in April

- Need to address also Tier2s
  - Cannot be done in the same way!
  - Work through existing structures, e.g.
  - HEPiX, national and regional bodies etc.
    - e.g. GridPP (12)

- Talking of SC Update at May HEPiX (FZK) with more extensive programme at Fall HEPiX (SLAC)
  - Maybe some sort of North American T2-fest around this?

# LCG Service Challenges - Overview

- LHC will enter production (physics) in April 2007
    - Will generate an enormous volume of data
    - Will require huge amount of processing power

- LCG 'solution' is a world-wide Grid
    - Many components understood, deployed, tested..

- But...
    - Unprecedented scale
    - Humungous challenge of getting large numbers of institutes and individuals, all with existing, sometimes conflicting commitments, to work together

- LCG must be ready at full production capacity, functionality and reliability in less than 2 years from now
    - Issues include h/w acquisition, personnel hiring and training, vendor rollout schedules etc.

- **Should not limit ability of physicist to exploit performance of detectors nor LHC's physics potential**
    - Whilst being stable, reliable and easy to use

# Key Principles

- Service challenges results in a **series** of services that exist in **parallel** with **baseline production** service

- Rapidly and successively approach production needs of LHC

- Initial focus: core (data management) services

- Swiftly expand out to cover **full spectrum** of production and analysis chain

- Must be as realistic as possible, including end-end testing of key experiment **use-cases** over extended periods with recovery from **glitches** and **longer-term** outages

➢ **Necessary resources and commitment pre-requisite to success!**

- Should not be under-estimated!

# Initial Schedule (1/2)

- Tentatively suggest quarterly schedule with monthly reporting
  - e.g. Service Challenge Meetings / GDB respectively
  - **Less than 7 complete cycles to go!**

- Urgent to have detailed schedule for 2005 with at least an outline for remainder of period asap
  - e.g. end January 2005

- Must be developed together with key partners
  - Experiments, other groups in IT, T1s, …

- Will be regularly refined, ever increasing detail…

- Detail must be such that partners can develop their own internal plans and to say what is and what is not possible
  - e.g. FIO group, T1s, …

# Initial Schedule (2/2)

- Q1 / Q2: up to 5 T1s, writing to tape at 50MB/s per T1 (no expts)

- Q3 / Q4: include two experiments and a few selected T2s

- 2006: progressively add more T2s, more experiments, ramp up to twice nominal data rate

- 2006: production usage by all experiments at reduced rates (cosmics); validation of computing models

- 2007: delivery and contingency

➢ N.B. there is more detail in December GDB presentations

➢ Need to be re-worked now!

# Agenda

_LCG Project, Grid Deployment Group, CERN_

- Reminder of the Goals and Timelines of the LCG Service Challenges

➢ **<u>Summary of LHC Experiments' Computing Models</u>**

- Outline of Service Challenges

- Review of SC1

- Status of SC2

- Plans for SC3 and beyond

# Computing Model Summary - Goals

- Present key features of LHC experiments' Computing Models in a consistent manner

- High-light the commonality

- Emphasize the key differences

- Define these 'parameters' in a central place (LCG web)
  - Update with change-log as required

- Use these parameters as input to requirements for Service Challenges

- To enable partners (T0/T1 sites, experiments, network providers) to have a clear understanding of what is required of them

- <u>Define precise terms and 'factors'</u>
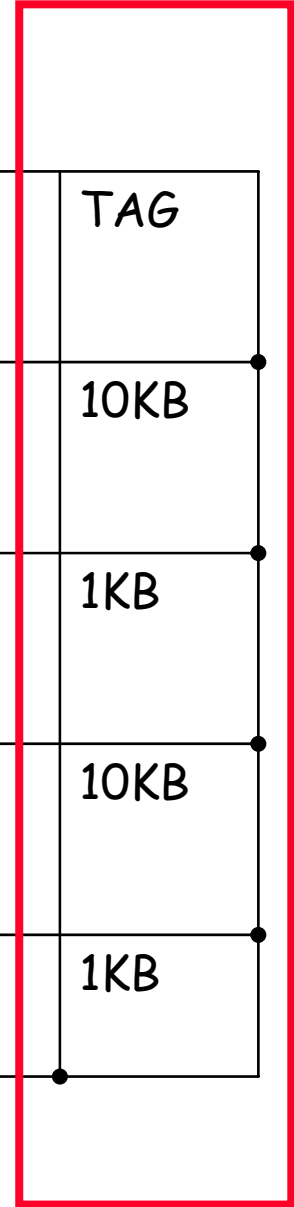
# Where do these numbers come from?

- Based on Computing Model presentations given to GDB in December 2004 and to T0/T1 networking meeting in January 2005
- Documents are those publicly available for January LHCC review
  - ☠ Official website is protected

- **Some details may change but the overall conclusions do not!**

- Part of plan is to understand how sensitive overall model is to variations in key parameters
- Iteration with experiments is on-going
  - i.e. I have tried to clarify any questions that I have had

- **Any mis-representation or mis-interpretation is entirely my responsibility**

- Sanity check: compare with numbers from MoU Task Force

| *LCG Project, Grid Deployment Group, CERN* | Nominal | These are the raw figures produced by multiplying e.g. event size x trigger rate. |
| | Headroom | A factor of 1.5 that is applied to cater for peak rates. |
| | Efficiency | A factor of 2 to ensure networks run at less than 50% load. |
| | Recovery | A factor of 2 to ensure that backlogs can be cleared within 24 – 48 hours and to allow the load from a failed Tier1 to be switched over to others. |
| | **Total Requirement** | **A factor of 6 must be applied to the nominal values to obtain the bandwidth that must be provisioned.**<br><br>**Arguably this is an over-estimate, as "Recovery" and "Peak load" conditions are presumably relatively infrequent, and can also be smoothed out using appropriately sized transfer buffers.**<br><br>But as there may be under-estimates elsewhere… |

All numbers presented will be
**nominal** unless explicitly specified

# Overview of pp running

| Experiment | SIM | SIMESD | RAW | Trigger | RECO | AOD | TAG |
|---|---|---|---|---|---|---|---|
| ALICE | 400KB | 40KB | 1MB | 100Hz | 200KB | 50KB | 10KB |
| ATLAS | 2MB | 500KB | 1.6MB | 200Hz | 500KB | 100KB | 1KB |
| CMS | 2MB | 400KB | 1.5MB | 150Hz | 250KB | 50KB | 10KB |
| LHCb | | 400KB | 25KB | 2KHz | 75KB | 25KB | 1KB |

# pp questions / uncertainties

- Trigger rates essentially independent of luminosity
  - Explicitly stated in both ATLAS and CMS CM docs

- Uncertainty (at least in my mind) on issues such as zero suppression, compaction etc of raw data sizes
  - Discussion of these factors in CMS CM doc p22:

- RAW data size ~300kB (Estimated from MC)
  - Multiplicative factors drawn from CDF experience
    - MC Underestimation factor 1.6
    - HLT Inflation of RAW Data, factor 1.25
    - Startup, thresholds, zero suppression,…. Factor 2.5
  - Real initial event size more like **1.5MB**
    - Could be anywhere between 1 and 2 MB
  - Hard to deduce when the even size will fall and how that will be compensated by increasing Luminosity

- i.e. total factor = 5 for CMS raw data

- N.B. must consider not only Data Type (e.g. result of Reconstruction) but also how it is used
  - e.g. compare how Data Types are used in LHCb compared to CMS

- All this must be plugged into the meta-model!

# Overview of Heavy Ion running

| Experiment | SIM | SIMESD | RAW | Trigger | RECO | AOD | TAG |
|---|---|---|---|---|---|---|---|
| ALICE | 300MB | 2.1MB | 12.5MB | 100Hz | 2.5MB | 250KB | 10KB |
| ATLAS | | | 5MB | 50Hz | | | |
| CMS | | | 7MB | 50Hz | 1MB | 200KB | TBD |
| LHCb | N/A | N/A | N/A | N/A | N/A | N/A | N/A |

# Heavy Ion Questions / Uncertainties

- Heavy Ion computing models less well established than for pp running

- I am concerned about model for 1st/2nd/3rd pass reconstruction and data distribution

➢ *"We therefore require that these data (Pb-Pb) are reconstructed at the CERN T0 and exported over a four-month period after data taking. This should leave enough time for a second and third reconstruction pass at the Tier 1's"* **(ALICE)**

- Heavy Ion model has major impact on those Tier1's supporting these experiments
  - All bar LHCb!

- Critical to clarify these issues as soon as possible…

# Data Rates from MoU Task Force

| MB/Sec | RAL | FNAL | BNL | FZK | IN2P3 | CNAF | PIC | T0 Total |
|---|---|---|---|---|---|---|---|---|
| ATLAS | 106.87 | 0.00 | 173.53 | 106.87 | 106.87 | 106.87 | 106.87 | 707.87 |
| CMS | 69.29 | 69.29 | 0.00 | 69.29 | 69.29 | 69.29 | 69.29 | 415.71 |
| ALICE | 0.00 | 0.00 | 0.00 | 135.21 | 135.21 | 135.21 | 0.00 | 405.63 |
| LHCb | 6.33 | 0.00 | 0.00 | 6.33 | 6.33 | 6.33 | 6.33 | 31.67 |
| T1 Totals MB/sec | 182.49 | 69.29 | 173.53 | 317.69 | 317.69 | 317.69 | 182.49 | 1560.87 |
| T1 Totals Gb/sec | 1.46 | 0.55 | 1.39 | 2.54 | 2.54 | 2.54 | 1.46 | 12.49 |
| | | | | | | | | |
| | | | | | | | | |
| Estimated T1 Bandwidth Needed | | | | | | | | |
| (Totals * 1.5*(headroom))*2(capacity)* | 4.38 | 1.66 | 4.16 | 7.62 | 7.62 | 7.62 | 4.38 | 37.46 |
| | | | | | | | | |
| | | | | | | | | |
| **Assumed Bandwidth Provisioned** | **10.00** | **10.00** | **10.00** | **10.00** | **10.00** | **10.00** | **10.00** | **70.00** |

*Spreadsheet used to do this calculation will be on Web.*

*Table is in*

*http://cern.ch/LCG/MoU%20meeting%20March%2010/Report_to_the_MoU_Task_Force.doc*

# Data Rates using CM Numbers

Steps:

- Take Excel file used to calculate MoU numbers

- Change one by one the Data Sizes as per latest CM docs

- See how overall network requirements change

➢ Need also to confirm that model correctly reflects latest thinking

➢ And understand how sensitive the calculations are to e.g. changes in RAW event size, # of Tier1s, roles of specific Tier1s etc.

# Base Requirements for T1s

- ➢ **Provisioned bandwidth comes in units of 10Gbits/sec although this is an evolving parameter**

  - ▪ *From* Reply to Questions from Computing MoU Task Force…

- ▪ Since then, some parameters of the Computing Models have changed

- ▪ Given the above quantisation, relatively insensitive to small-ish changes

- ▪ Important to understand implications of multiple-10Gbit links, particularly for sites with Heavy Ion programme

- ➢ **<u>For now, need plan for 10Gbit links to all Tier1s</u>**

# Response from 'Networkers'

[Hans Döbbeling] believe the GEANT2 consortium will be able to deliver the following for the 7 European TIER1s:

1.  list of networking domains and technical contacts
2.  time plan of  availability of services 1G VPN, 1G Lambda, 10Gig Lambda at the national GEANT2 pops and at the TIER1 sites.
3.  a model for SLAs and the monitoring of SLAs
4.  a proposal for operational procedures
5.  a compilation of possible cost sharing per NREN

Proposes that CERN focuses on issues related to non-European T1s

# Agenda

*LCG Project, Grid Deployment Group, CERN*

- Reminder of the Goals and Timelines of the LCG Service Challenges

- Summary of LHC Experiments' Computing Models

- Outline of Service Challenges

➢ **Review of SC1**

- Status of SC2

- Plans for SC3 and beyond

# "Review of Service Challenge 1"

James Casey, IT-GD, CERN

RAL, 26 January 2005

# Overview

- Reminder of targets for the Service Challenge

- What we did

- What can we learn for SC2?

# Milestone I & II Proposal

From NIKHEF/SARA Service Challenge Meeting

*LCG Project, Grid Deployment Group, CERN*

# Service Challenge Schedule

LCG Project,  Grid Deployment Group, CERN

From FZK Dec Service Challenge Meeting:

# SARA – Dec 04

- **Used a SARA specific solution**
  - Gridftp running on 32 nodes of SGI supercomputer (teras)
  - 3 x 1Gb network links direct to teras.sara.nl
  - 3 gridftp servers, one for each link
  - Did load balancing from CERN side
    - 3 oplapro machines transmitted down each 1Gb link
  - Used radiant-load-generator script to generate data transfers
- **Much efforts was put in from SARA personnel (~1-2 FTEs) before and during the challenge period**
- **Tests ran from 6-20<sup>th</sup> December**
  - Much time spent debugging components

# Problems seen during SC1

- **Network Instability**
  - Router electrical problem at CERN
  - Interruptions due to network upgrades on CERN test LAN
- **Hardware Instability**
  - Crashes seen on teras 32-node partition used for challenges
  - Disk failure on CERN transfer node
- **Software Instability**
  - Failed transfers from gridftp. Long timeouts resulted in significant reduction in throughput
  - Problems in gridftp with corrupted files
- **Often hard to isolate a problem to the right subsystem**

# SARA SC1 Summary

- Sustained run of 3 days at end
  - 6 hosts at CERN side. single stream transfers, 12 files at a time
  - Average throughput was 54MB/s
  - Error rate on transfers was 2.7%
- Could transfer down each individual network links at ~40MB/s
  - This did not translate into the expected 120MB/s speed
  - Load on teras and oplapro machines was never high (~6-7 for a 32 node teras, < 2 for 2-node oplapro) Load on oplapro machines
- See Service Challenge wiki for logbook kept during Challenge

# Gridftp problems

- **64 bit compatibility problems**
  - logs negative numbers for file size > 2 GB
  - logs erroneous buffer sizes to the logfile if the server is 64-bits
- **No checking of file length on transfer**
  - No error message doing a third party transfer with corrupted files
- **Issues followed up with globus gridftp team**
  - First two will be fixed in next version.
  - The issue of how to signal problems during transfers is logged as an enhancement request

# FermiLab & FZK – Dec 04/Jan 05

- **FermiLab declined to take part in Dec 04 sustained challenge**
    - They had already demonstrated 500MB/s for 3 days in November



- **FZK started this week**
    - Bruno Hoeft will give more details in his site report

# What can we learn ?

- ☠ <u>**SC1 did not succeed**</u>
  - ▪ We did not meet the milestone of 500MB/s for 2 weeks
- ▪ We need to do these challenges to see what actually goes wrong
  - ▪ A lot of things do, and did, go wrong
- ▪ We need better test plans for validating the infrastrcture before the challenges (network throughput, disk speeds, etc…)
  - ▪ Ron Trompert (SARA) has made a first version of this
- ▪ We need to proactively fix low-level components
  - ▪ Gridftp, etc…
- ➢ <u>**SC2 and SC3 will be a lot of work !**</u>

# 2005 Q1(i)

SC2 - Robust Data Transfer Challenge

Set up infrastructure for 6 sites
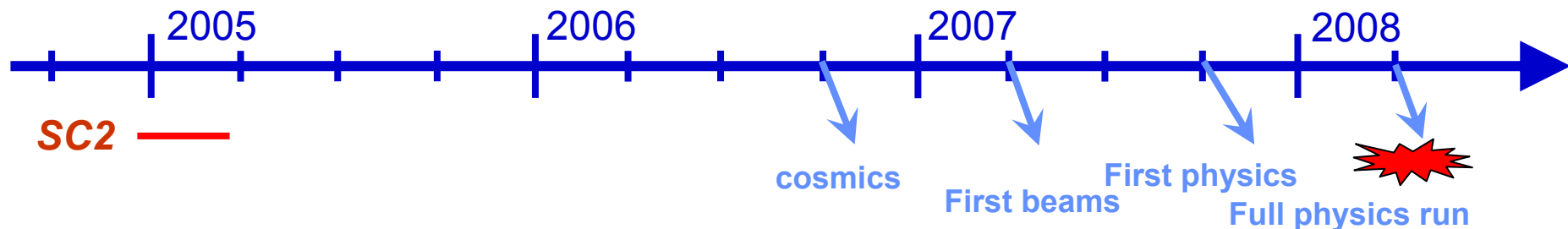- Fermi, NIKHEF/SARA, GridKa, RAL, CNAF, CCIN2P3

Test sites individually
 – at least two at 500 MByte/s with CERN

Agree on sustained data rates for each participating centre
Goal – by end March sustained 500 Mbytes/s aggregate at CERN

In parallel - serve the ATLAS "Tier0 tests" (needs more discussion)

*LCG Project, Grid Deployment Group, CERN*

2005    2006    2007    2008

SC2 ——

cosmics

First beams

First physics

Full physics run

# Status of SC2

# 2005 Q1(ii)

RADIANT

In parallel with SC2
– prepare for the next service challenge (SC3)

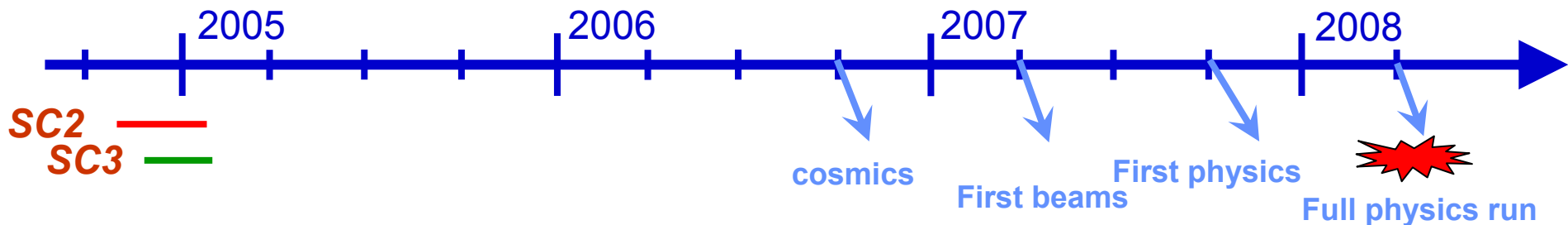Build up 1 GByte/s *challenge* facility at CERN
- **The current 500 MByte/s facility used for SC2 will become the *testbed* from April onwards (10 ftp servers, 10 disk servers, network equipment)**

Build up infrastructure at each external centre
- **Average capability ~150 MB/sec at a Tier-1 (to be agreed with each T-1)**

Further develop reliable transfer framework software
- **Include catalogues, include VO's**

2005        2006        2007        2008

SC2 —
SC3 —

cosmics

First beams

First physics

Full physics run

# 2005 Q2-3(i)

## SC3 - 50% service infrastructure

- Same T1s as in SC2 (Fermi, NIKHEF/SARA, GridKa, RAL, CNAF, CCIN2P3)
- Add at least two T2s
- "50%" means approximately 50% of the nominal rate of ATLAS+CMS
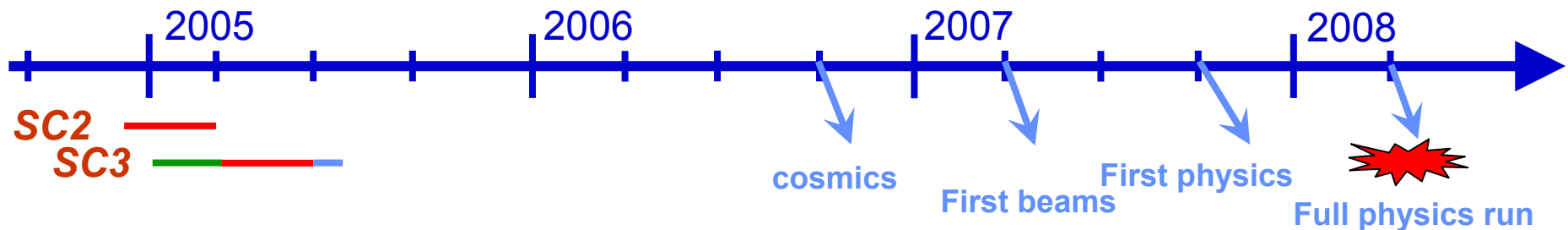
Using the 1 GByte/s *challenge* facility at CERN -
- Disk at T0 to tape at all T1 sites at 80 Mbyte/s
- Data recording at T0 from same disk buffers
- Moderate traffic disk-disk between T1s and T2s

**Use ATLAS and CMS files, reconstruction, ESD skimming codes**

Goal - 1 month sustained service in July
- 500 MBytes/s aggregate at CERN, 80 MBytes/s at each T1

2005        2006        2007        2008

SC2

SC3

cosmics

First beams

First physics

Full physics run

# SC3 Planning

- Meetings with other IT groups at CERN to refine goals of SC3 (milestone document) and steps that are necessary to reach them

- IT-GM: definition of middleware required, schedule, acceptance tests etc

- IT-ADC: pre-production and production services required, e.g. Database backends

- IT-FIO: file transfer servers etc etc

- IT-CS: network requirements etc

- Informal discussions with experiments regarding involvement of production teams

- Many details missing: (one being identification of T2s)

# SC3 Planning - cont

- Base (i.e. non-experiment) software required for SC3 scheduled for delivery end-February 2005

- Targeting service infrastructure for same date
  - Database services, Gridftp servers etc.

- Acceptance testing during March

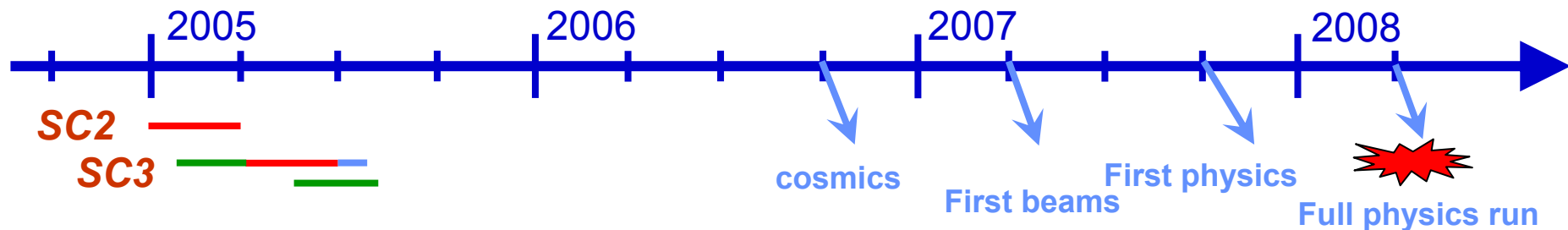- In parallel, have to start discussions with experiments, T1s etc

# 2005 Q2-3(ii)

In parallel with SC3 prepare additional centres using the 500 MByte/s test facility

- **Test Taipei, Vancouver, Brookhaven, additional Tier-2s**

Further develop framework software

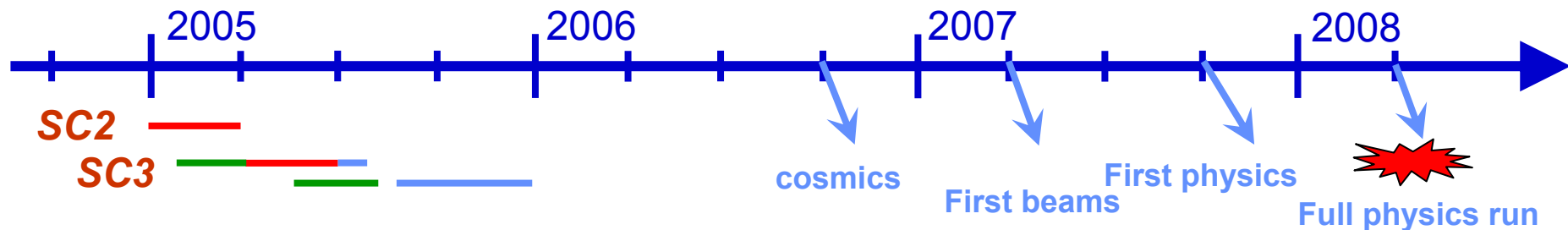- **Catalogues, VO's, use experiment specific solutions**

2005          2006          2007          2008

**SC2**

**SC3**

cosmics

**First beams**

**First physics**

**Full physics run**

# 2005 – September-December (i)

**50% Computing Model Validation Period**

The service exercised in SC3 is made available to experiments for computing model tests

Additional sites are added as they come up to speed

End-to-end data rates –
- 500 Mbytes/s at CERN (aggregate)
- 80 Mbytes/s at Tier-1s
- Modest Tier-2 traffic
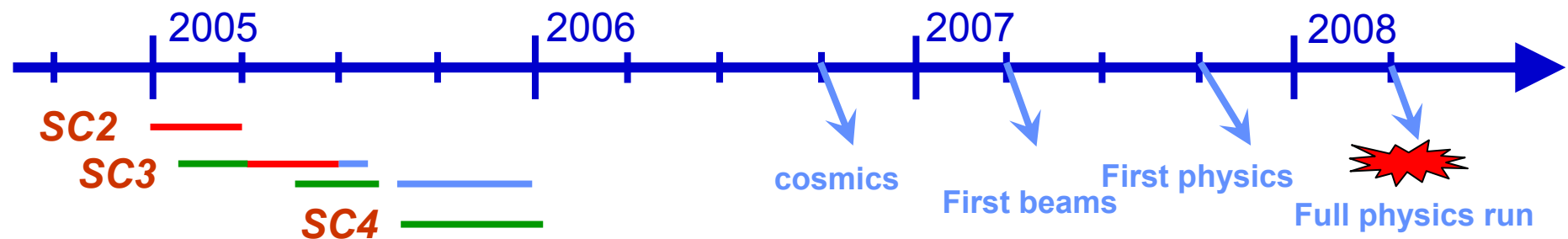
# 2005 – September-December (ii)

In parallel with the SC3 model validation period,
in preparation for the first 2006 service challenge (SC4) –

Using 500 MByte/s test facility
- **test PIC and Nordic T1s**
- **and T2's that are ready (Prague, LAL, UK, INFN, ..**

Build up the production facility at CERN to 3.6 GBytes/s

Expand the capability at all Tier-1s to full nominal data rate



2005          2006          2007          2008

SC2

SC3

SC4

cosmics

First beams

First physics

Full physics run

# 2006 - January-August
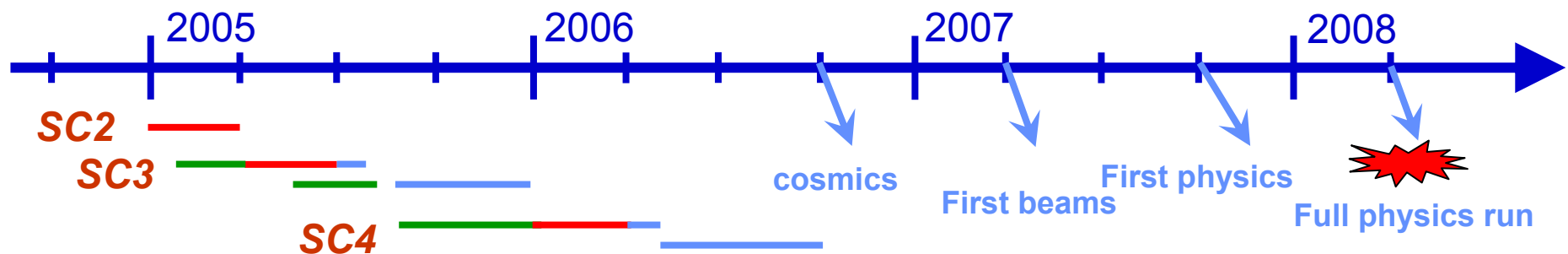
RADIANT

## SC4 – full computing model services

- Tier-0, ALL Tier-1s, all major Tier-2s operational
  at full target data rates (~1.8 GB/sec at Tier-0)
- acquisition - reconstruction - recording – distribution,
  **PLUS** ESD skimming, servicing Tier-2s

Goal – stable test service for one month – April 2006

## 100% Computing Model Validation Period (May-August 2006)

Tier-0/1/2 full model test - **All experiments**

- 100% nominal data rate, with processing load scaled to 2006 cpus

2005      2006      2007      2008

**SC2**

**SC3**

**SC4**

cosmics

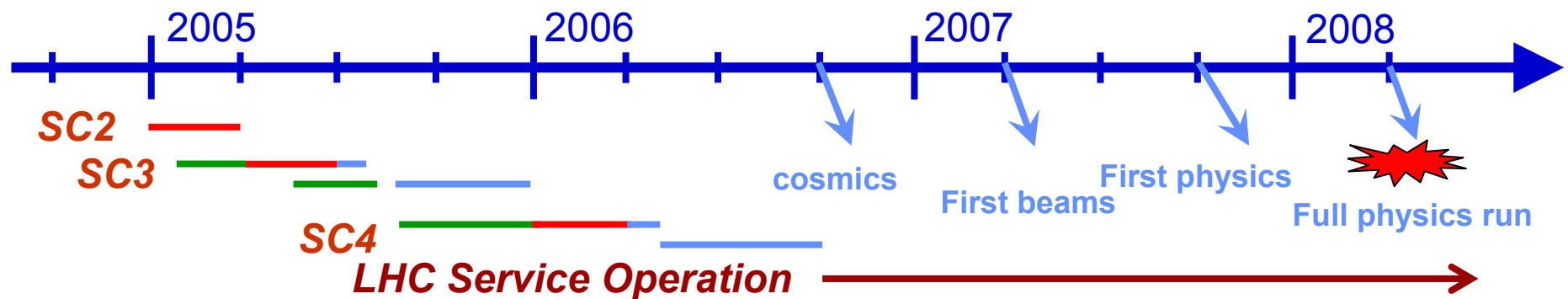First beams

First physics

Full physics run

# 2006 - September

The SC4 service becomes the permanent LHC service – available for experiments' testing, commissioning, processing of cosmic data, etc.

All centres ramp-up to capacity needed at LHC startup
- TWICE nominal performance
- Milestone to demonstrate this 6 months before first physics data

# Timeline

- Official target date for first collisions in LHC: April 2007

- Including ski-week(s), this is only 2 years away!

- But the real target is even earlier!

- Must be ready 6 months prior to data taking

- And data taking starts earlier than colliding beams!

- Cosmics (ATLAS in a few months), calibrations, single beams, …

# Conclusions

- To be ready to fully exploit LHC, significant resources need to be allocated to a series of service challenges by all concerned parties

- These challenges should be seen as an essential on-going and long-term commitment to achieving production LCG

- The countdown has started – we are already in (pre-)production mode

- Next stop: 2020