# CMS FacilitiesOps and IN2P3

## [ CMS visit to IN2P3 – Lyon, 23 Oct 09]

Daniele Bonacorsi, Peter Kreuzer
[ CMS Facilities Ops ]

Claudio Grandi, Chris Brew
[ T1 coordination in CMS Facilities Ops ]

Andrea Sciabà, Josep Flix
[ Site Readiness in CMS Facilities Ops ]

Nicolò Magini
[ Data Transfer Operations and DDT in CMS Facilities Ops ]

# CMS FacilitiesOps

## CMS FacilitiesOps weekly meetings

- ✦ To discuss status of T1 and T2 sites, and related items, over last 7 days
  - CMS contacts at T1's asked to provide brief weekly reports
  - SAM and SiteReadiness status is reviewed, explanations are asked, discussion
- ✦ Weekly, Monday afternoon, 5pm GVA time

## CMS attends WLCG Ops daily calls, 3pm GVA time

- ✦ Official WLCG official minutes:
  - https://twiki.cern.ch/twiki/bin/view/LCG/WLCGOperationsMeetings
- ✦ Collection of CMS daily reports:
  - https://twiki.cern.ch/twiki/bin/view/CMS/FacOps_WLCGdailyreports

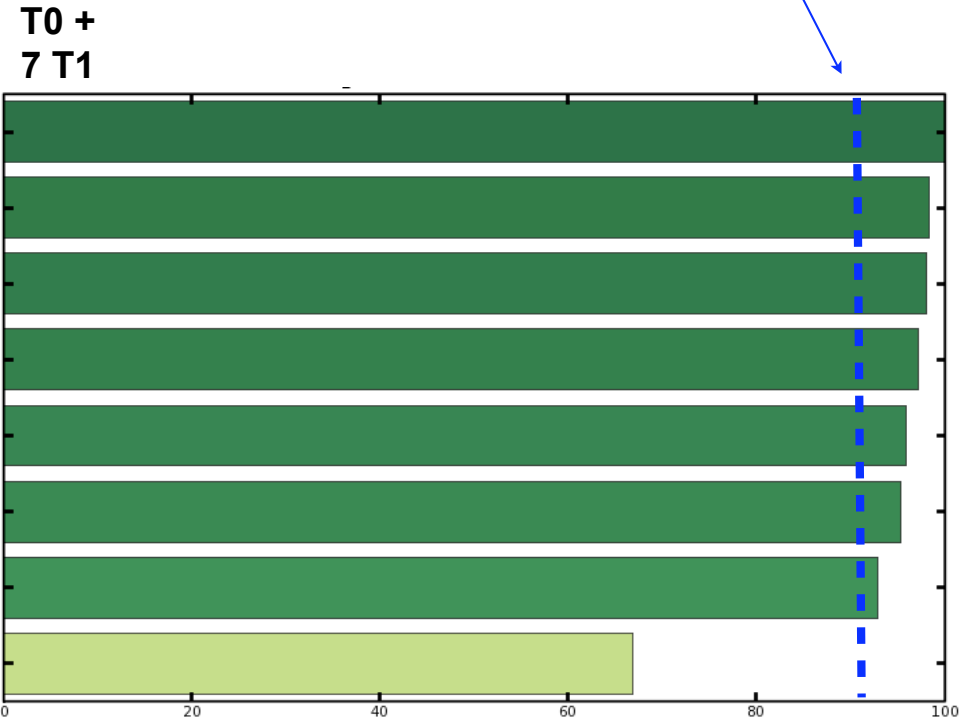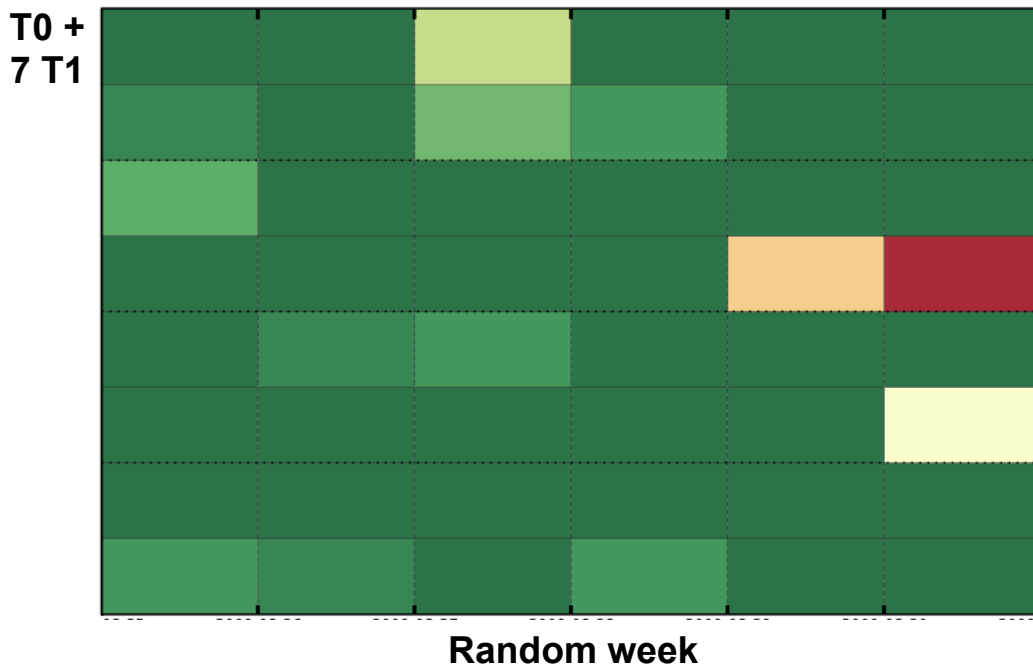# SAM Availability for CMS T1's

## CMS-specific SAM tests

✦ Complementary to WLCG SAM, to mimic real CMS workflows

- Widely documented elsewhere

## Overall SAM Availability ranking for CMS T1's: goal is **90%**

✦ For all orangish/redish boxes we discuss at FacOps weekly meetings



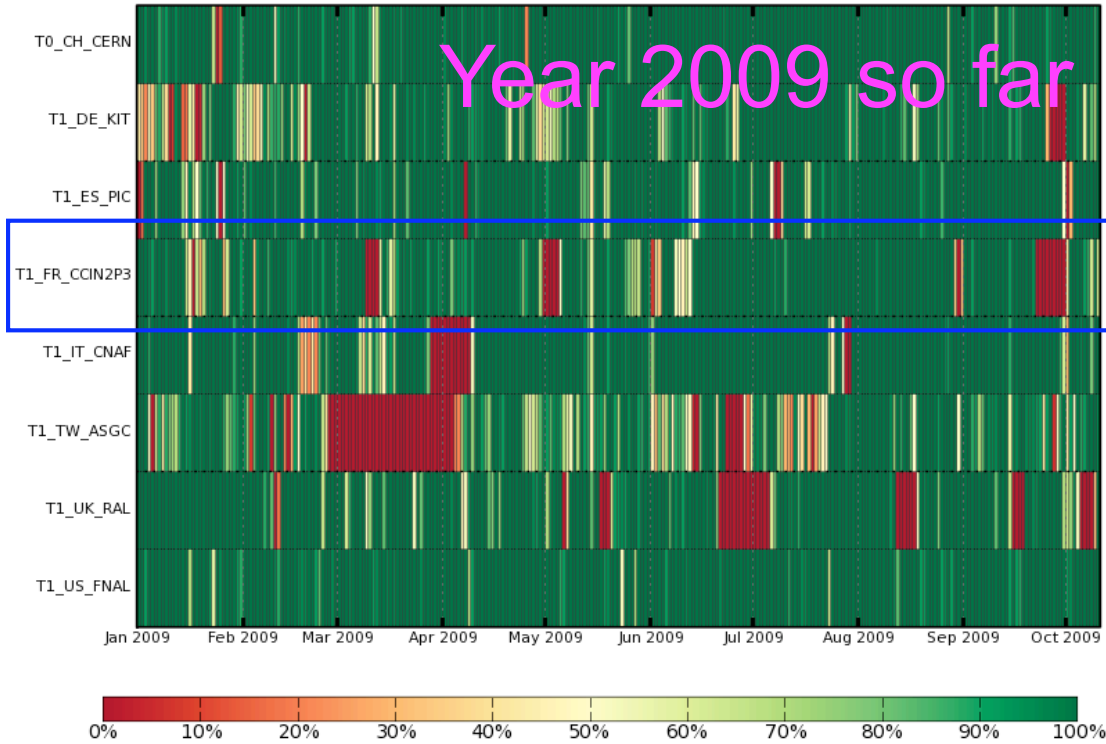T0 + 7 T1

**Random week**

T0 + 7 T1

# SAM Availability for *IN2P3*

## Looking to IN2P3 in CMS-specific SAM tests in 2009

# CMS SiteReadiness

## Global estimator in FacOps for the readiness of sites for daily operations



| | T1_[region]_[sitename] | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Site Readiness Status:** | R | W | R | W | NR | NR | NR | R | R | NR | NR | NR | NR | NR | R |
| **Daily Metric:** | O | O | O | O | O | O | O | O | E | O | E | E | E | O | O | O | E | E | O | E | O | O |
| **Maintenance:** | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up | Up |
| **Job Robot:** | 98% | 96% | 99% | 98% | 97% | 96% | 98% | 96% | 89% | 98% | 97% | 96% | 85% | 100% | 100% | 100% | 80% | 0% | 100% | 100% | 100% | 99% |
| **SAM Availability:** | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 88% | 60% | 28% | 100% | 100% | 100% | 96% | 100% | 100% | 80% | 100% | 100% |
| **T1::downlinkT0:** | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| **T1::downlinks/uplinksT1s:** | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) | 8(d)-8(u) |
| **T1::uplinksT2s:** | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 |
| | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 |
| | Sep | | | | | | | | | | | | Oct | | | | | | | | | |

https://twiki.cern.ch/twiki/bin/view/CMS/SiteCommRules

**"Site Readiness Status"** as defined in Site Commissioning Twiki:

- **R** = READY
- **W** = WARNING
- **NR** = NOT-READY
- **SD** = SCHEDULED-DOWNTIME

**"Daily Metric"** as boolean AND of all invidual metrics considered for the site:

- **O** = OK (All individual metrics above Site Commissioning Thresholds; "n/a" ignored)
- **E** = ERROR (Some individual metrics below Site Commissioning Thresholds)
- **SD** = SCHEDULED-DOWNTIME

- INDIVIDUAL METRICS -

**"Scheduled Downtimes":** site maintenances
- **Up** = Site is not declaring Scheduled-downtime
- = SD=full-site; SE-SD: All CMS SE(s) in SD; CE-SD: All CMS CE(s) in SD
- **~** = Some SE or CE services (not all) Downtime

**"SAM Availability":**
- = SAM availability is ≥ 90%
- = SAM availability is < 90%

**"T1::downlinks/uplinksT1s":**
- = Site has ≥ 4 DDT commissioned uplinks and downlinks, respectively, with other T1 sites
- = Otherwise

**"Job Robot":**
- = Job success rate is ≥ 90%
- = Job success rate is < 90%
- **-** = Jobs submitted but not finished
- **n/a** = Job success rate is n/a

**"T1::downlinkT0":**
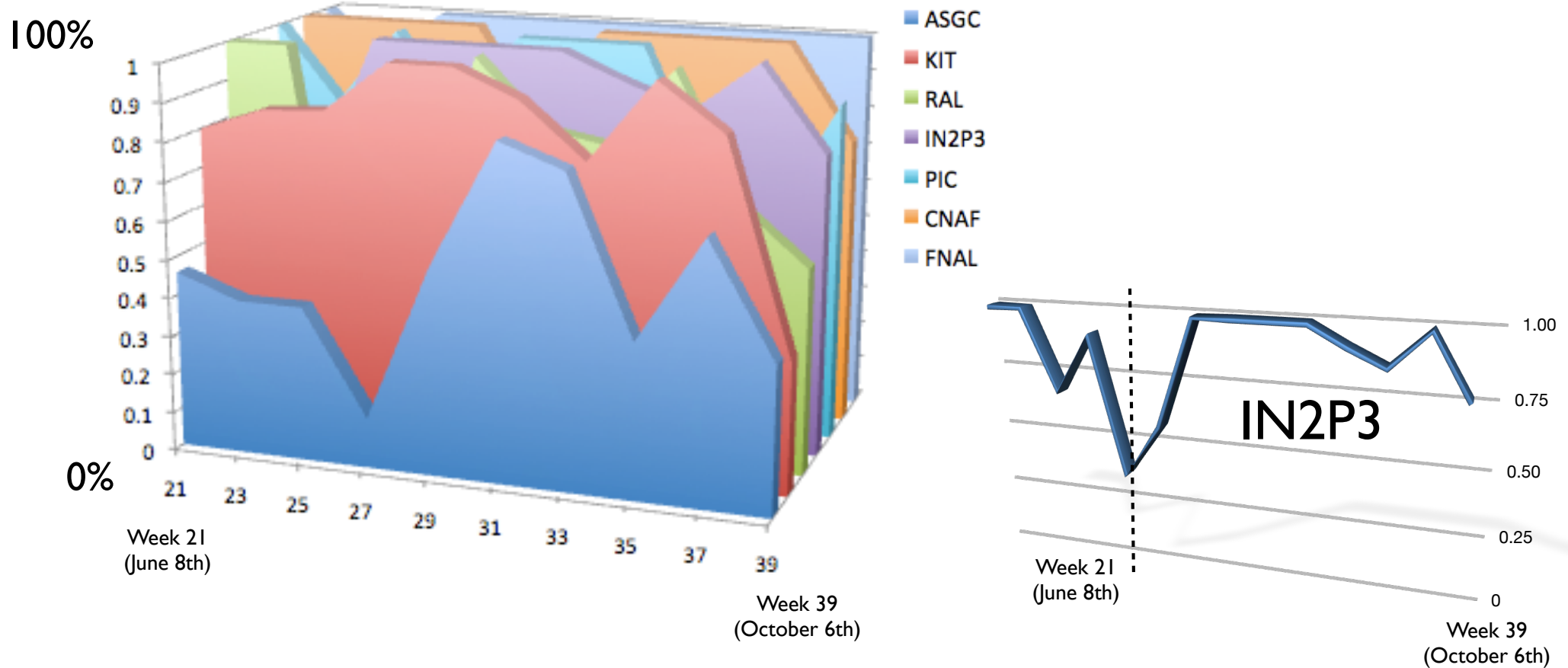- = Downlink from T0_CH_CERN is DDT-commissioned
- = Otherwise

**"T1::uplinksT2s":**
- = Site has ≥ 20 DDT commissioned uplinks to T2 sites
- = Otherwise

## From CMS Site Readiness metrics:

✦ Site availability: fraction of time all functional tests succeed

✦ JobRobot efficiency: fraction of successful "fake" analysis jobs

✦ Links: # of commissioned data transfer links

# CMS SiteReadiness ranking for CMS T1's



SiteReadiness goal for T1's: **90%**
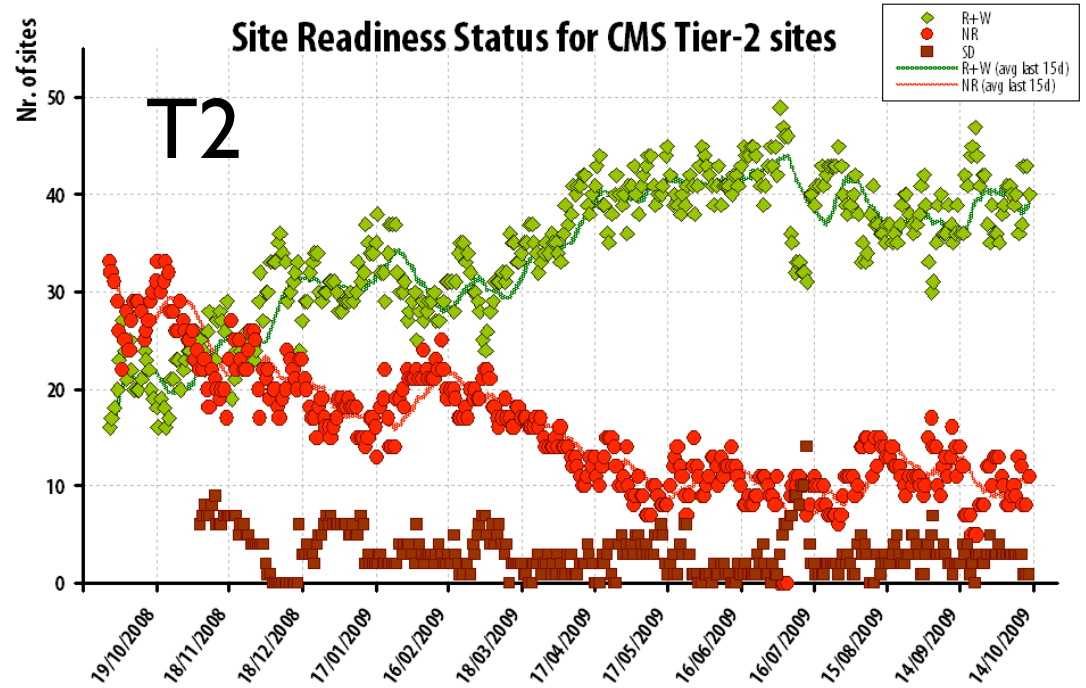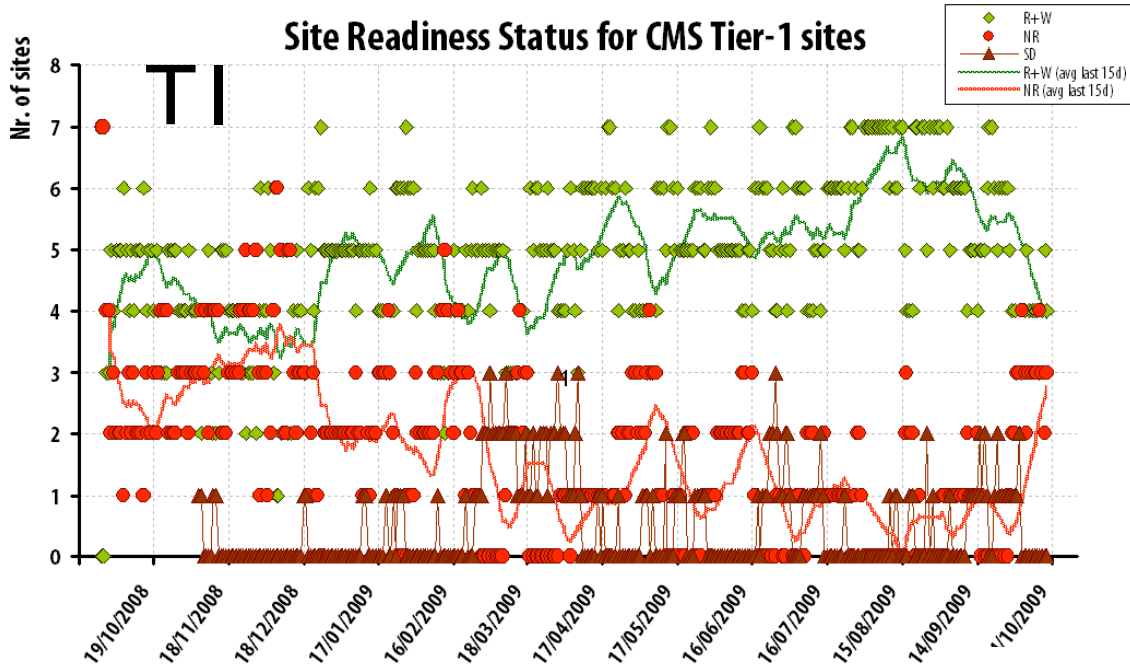
Achieved averages in *Jun-Oct 2009*:

- {FNAL, CNAF} at {**99%**, **95%**}
- {PIC, IN2P3, KIT, RAL} at {**87%**, **86%**, **85%**, **73%**}
- ASGC at **50%**

**WLCG SAM (ops)** not the full picture

**CMS-specific SAM** not the full picture

**SiteReadiness** (even!) not the full picture

- Need high CPU eff, disk stability, MSS solidity and performance, ...

*Example*: historical data on T1 and T2 sites

# CMS SiteReadiness ranking for *IN2P3*



Week 21
(June 8th)

Week 39
(October 6th)

# SiteReadiness breakdown for _IN2P3_

| Period / State | READY [days] | WARN [days] | NOT READY [days] | Downtime |
|---|---|---|---|---|
| June 09 | 20 | 0 | 10 | 0 |
| July 09 | 29 | 2 | 0 | 0 |
| Aug 09 | 27 | 3 | 0 | 1 |
| Sep 09 | 18 | 2 | 3 | 7 |

NOTE: SiteReadiness has lately suffered from SSB instabilities when tracing scheduled downtimes. The September IN2P3 downtime was corrected on SiteReadiness tables as announced here:

✦ https://hypernews.cern.ch/HyperNews/CMS/get/sc4/1969.html

# Transfer rates: T0 -> IN2P3



CMS PhEDEx - Transfer Rate
26 Weeks from Week 16 of 2009 to Week 42 of 2009

■ T0_CH_CERN_Export to T1_FR_CCIN2P3_Buffer

Maximum: 10.59 MB/s, Minimum: 0.00 MB/s, Average: 2.80 MB/s, Current: 1.12 MB/s

CMS PhEDEx - Transfer Quality
26 Weeks from Week 16 of 2009 to Week 42 of 2009

## Little activity in the PhEDEx /Prod instance

✦ few datasets from T0 assigned to IN2P3 as custodial
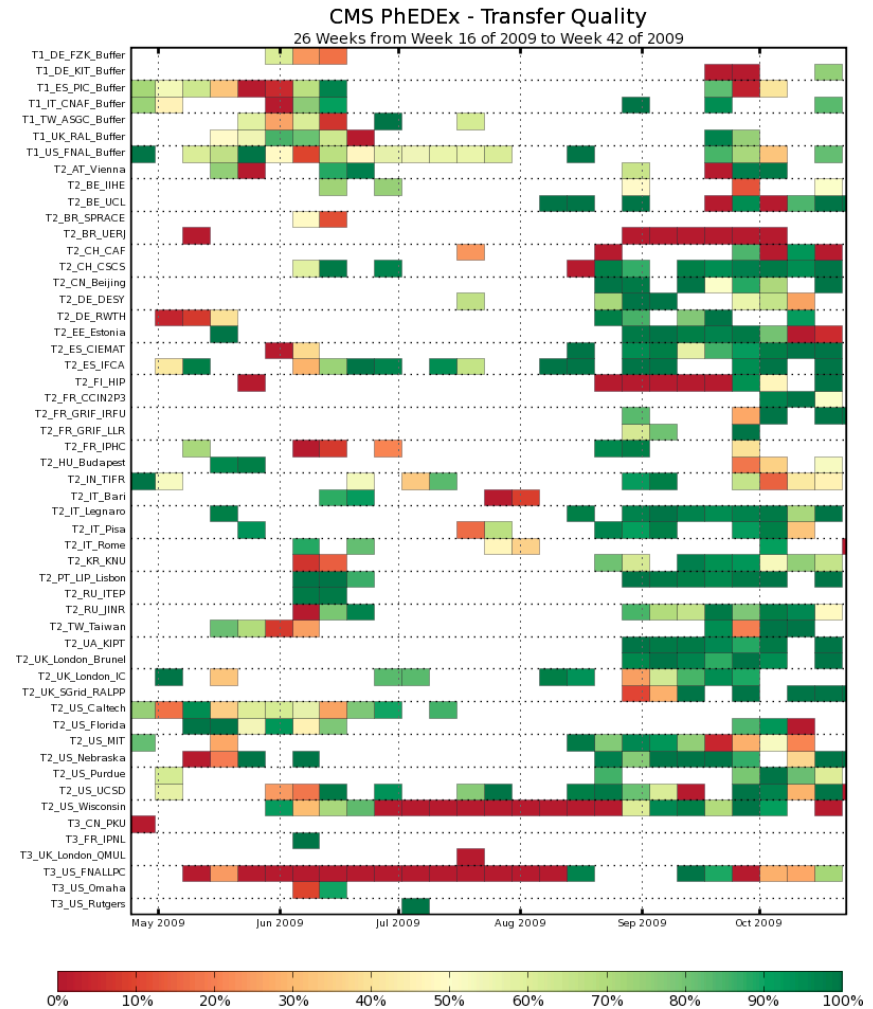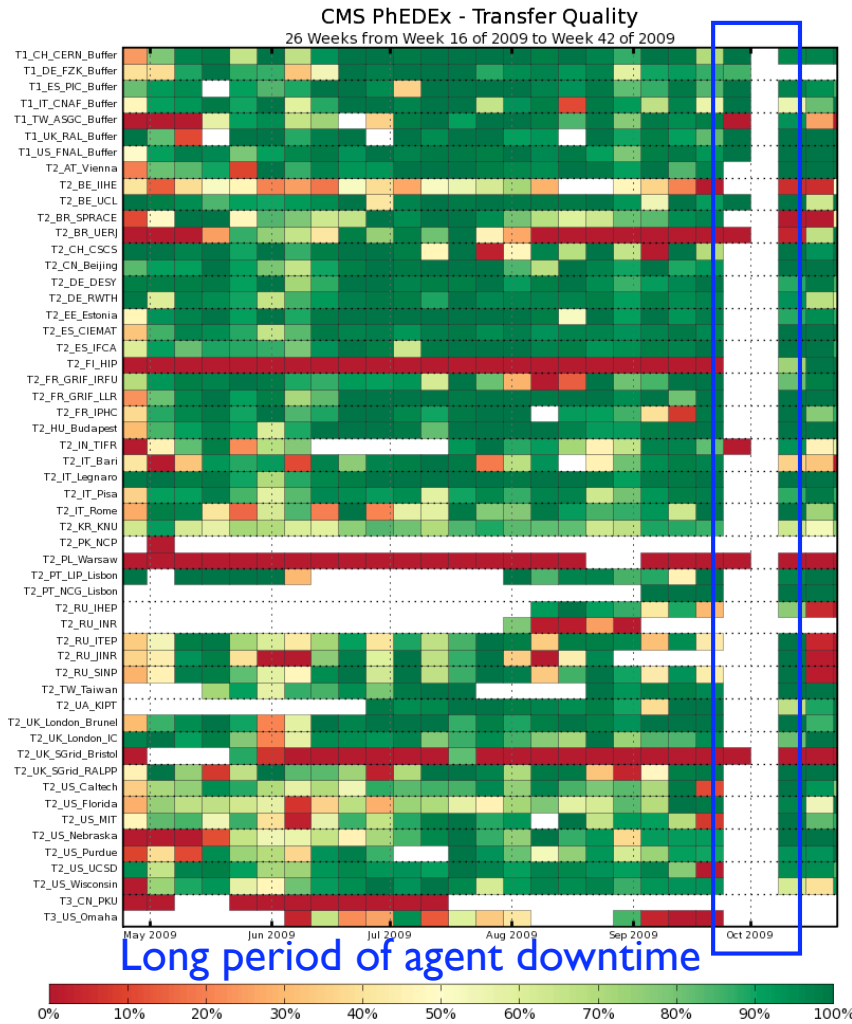site...

# Transfer rates: IN2P3 -> *



IN2P3 -> *
in **/Debug**

Long period of
agent downtime

[ Savannah #110535 ]

IN2P3 -> *
in **/Prod**

# Transfer quality: IN2P3 -> *

### IN2P3 -> * in **/Debug**



Long period of agent downtime
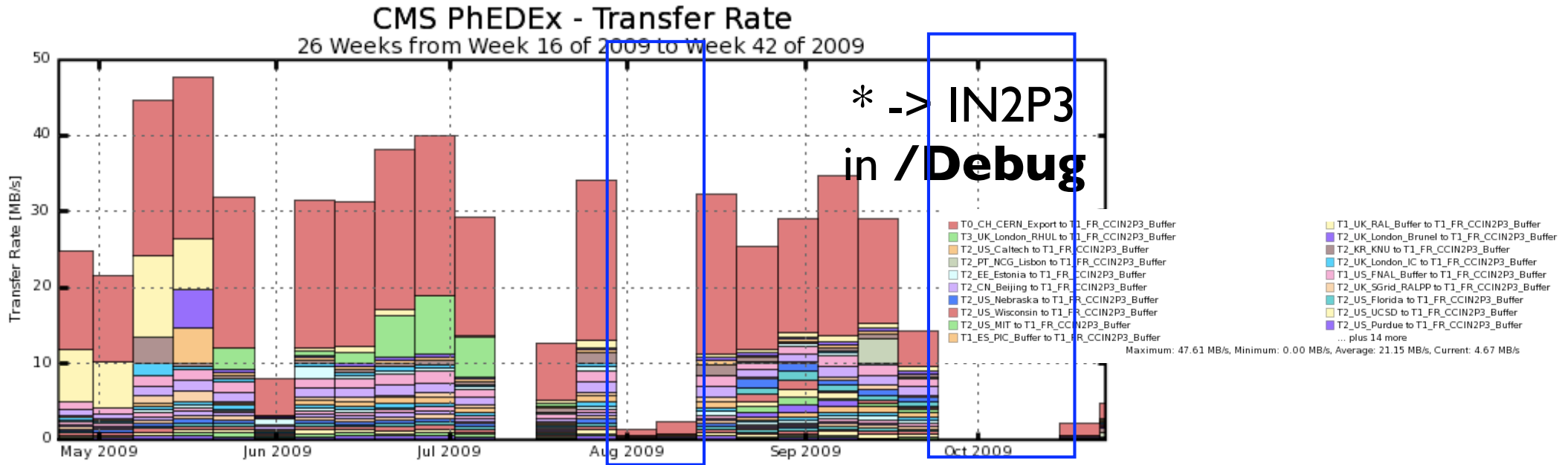
### IN2P3 -> * in **/Prod**



## Generally OK in /Debug

✦ apart from the agent downtime in late Sept

## Not too bad in /Prod

✦ Most frequent problem is transfer expirations due to FTS channel congestion - these are invisible in the plots...
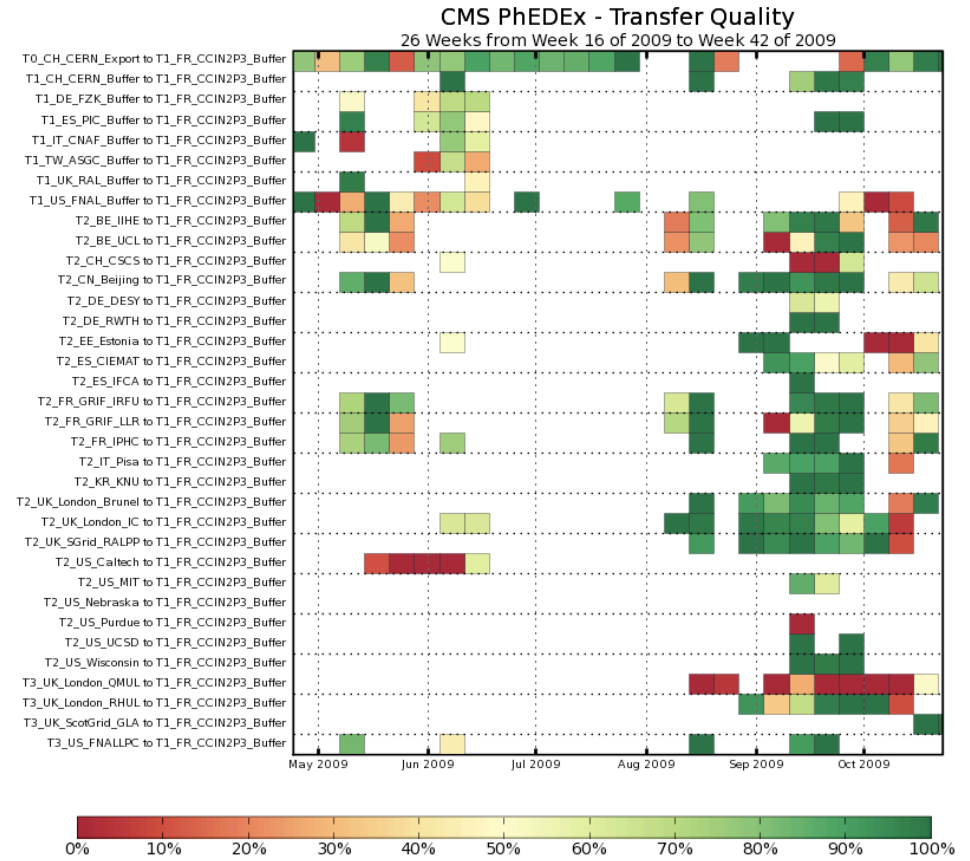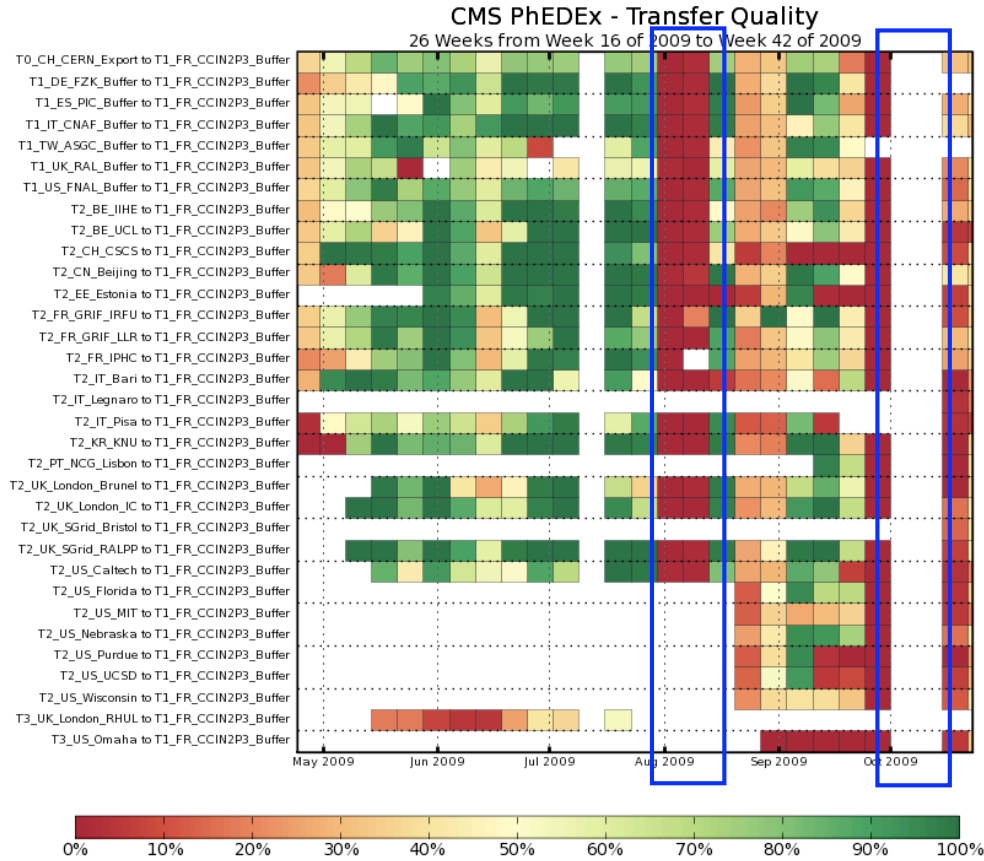
# Transfer rates: * -> IN2P3



* -> IN2P3
in **/Debug**

* -> IN2P3
in **/Prod**

# Transfer quality: * -> IN2P3

## * -> IN2P3 in /Debug



CMS PhEDEx - Transfer Quality
26 Weeks from Week 16 of 2009 to Week 42 of 2009

## * -> IN2P3 in /Prod



CMS PhEDEx - Transfer Quality
26 Weeks from Week 16 of 2009 to Week 42 of 2009

The import in the /Debug instance are more frequently in overall bad health

- ✦ agents down for long periods of time
- ✦ Relatively bad transfer quality in imports since summer

A large source of errors is "*Already have 1 record(s) with pnfsPath=[...]"

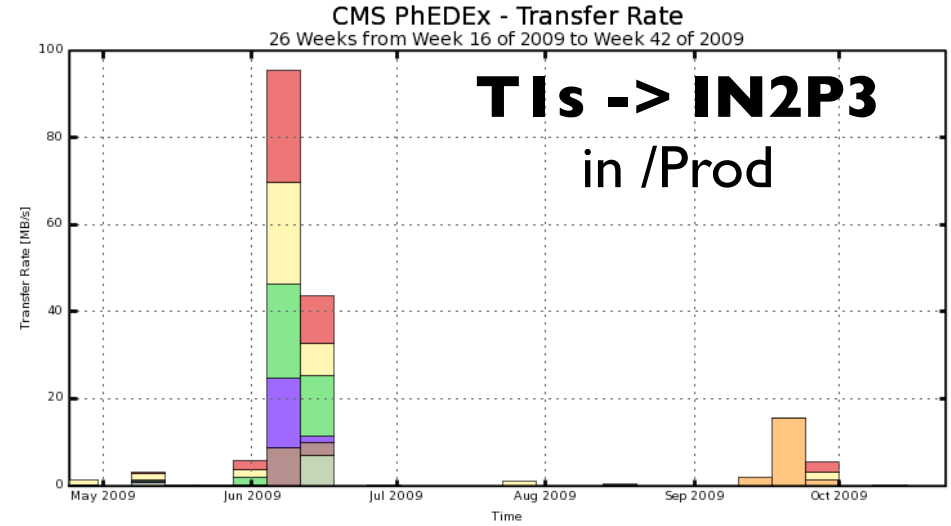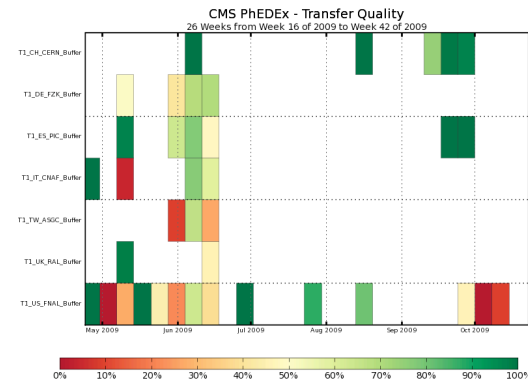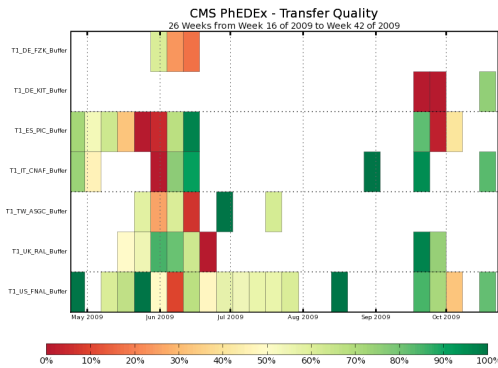- ✦ probably a cleanup of the LoadTest target area would improve things...
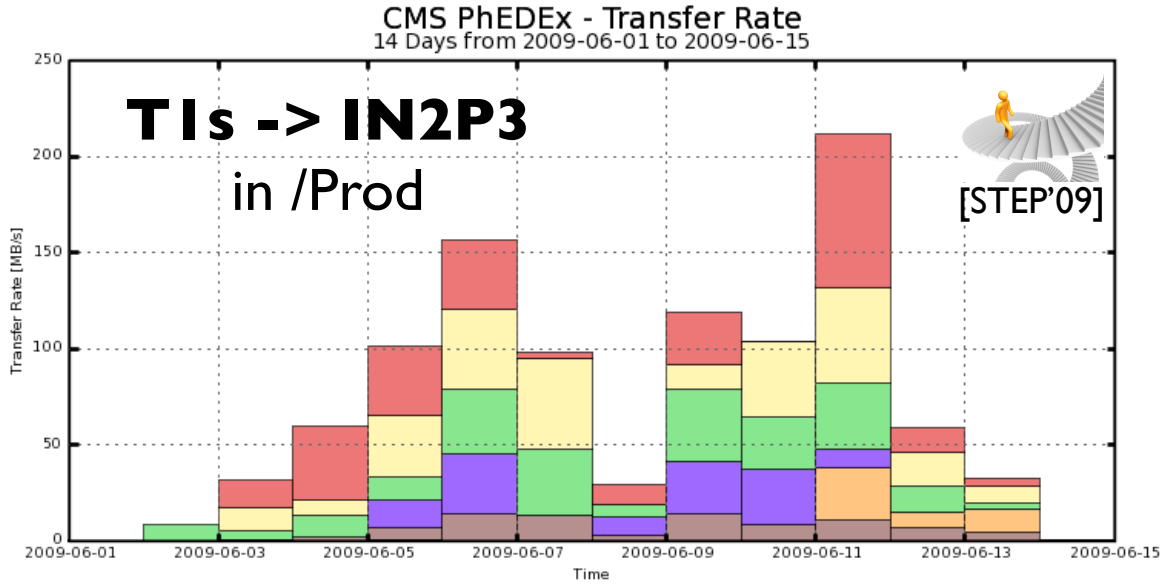
**Almost no activity outside STEP09**

**During STEP09, very good rates** (*more in the back-up slides*)

✦ Targets (assuming no rerouting in PhEDEx) were 185 MB/s in, 105 MB/s out - exceeded in one day
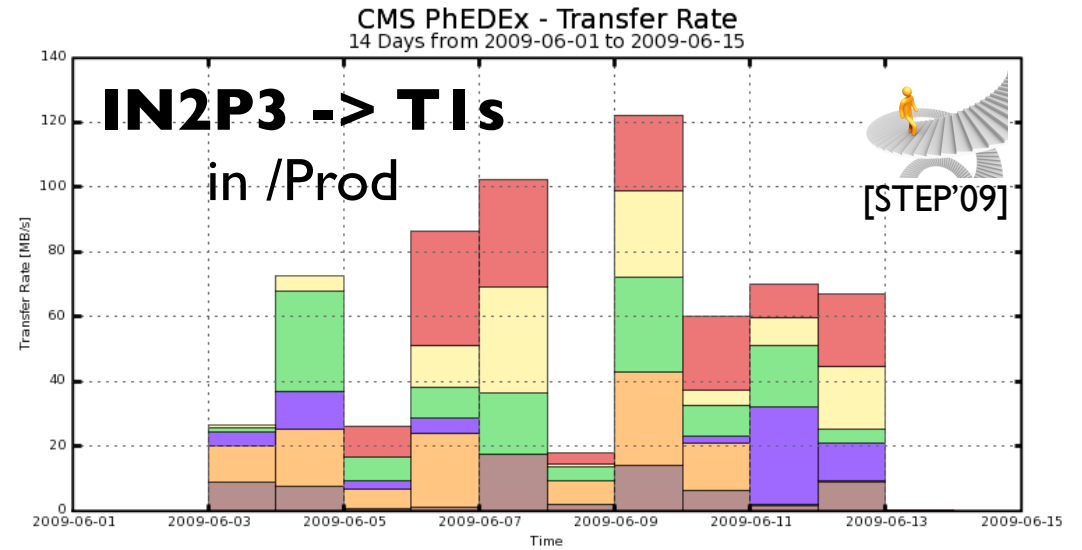
CMS PhEDEx - Transfer Rate
14 Days from 2009-06-01 to 2009-06-15

**T1s -> IN2P3** in /Prod

[STEP'09]

T1_DE_FZK_Buffer to T1_FR_CCIN2P3_Buffer
T1_US_FNAL_Buffer to T1_FR_CCIN2P3_Buffer
T1_UK_RAL_Buffer to T1_FR_CCIN2P3_Buffer
T1_CH_CERN_Buffer to T1_FR_CCIN2P3_Buffer
T1_ES_PIC_Buffer to T1_FR_CCIN2P3_Buffer
T1_IT_CNAF_Buffer to T1_FR_CCIN2P3_Buffer
T1_TW_ASGC_Buffer to T1_FR_CCIN2P3_Buffer

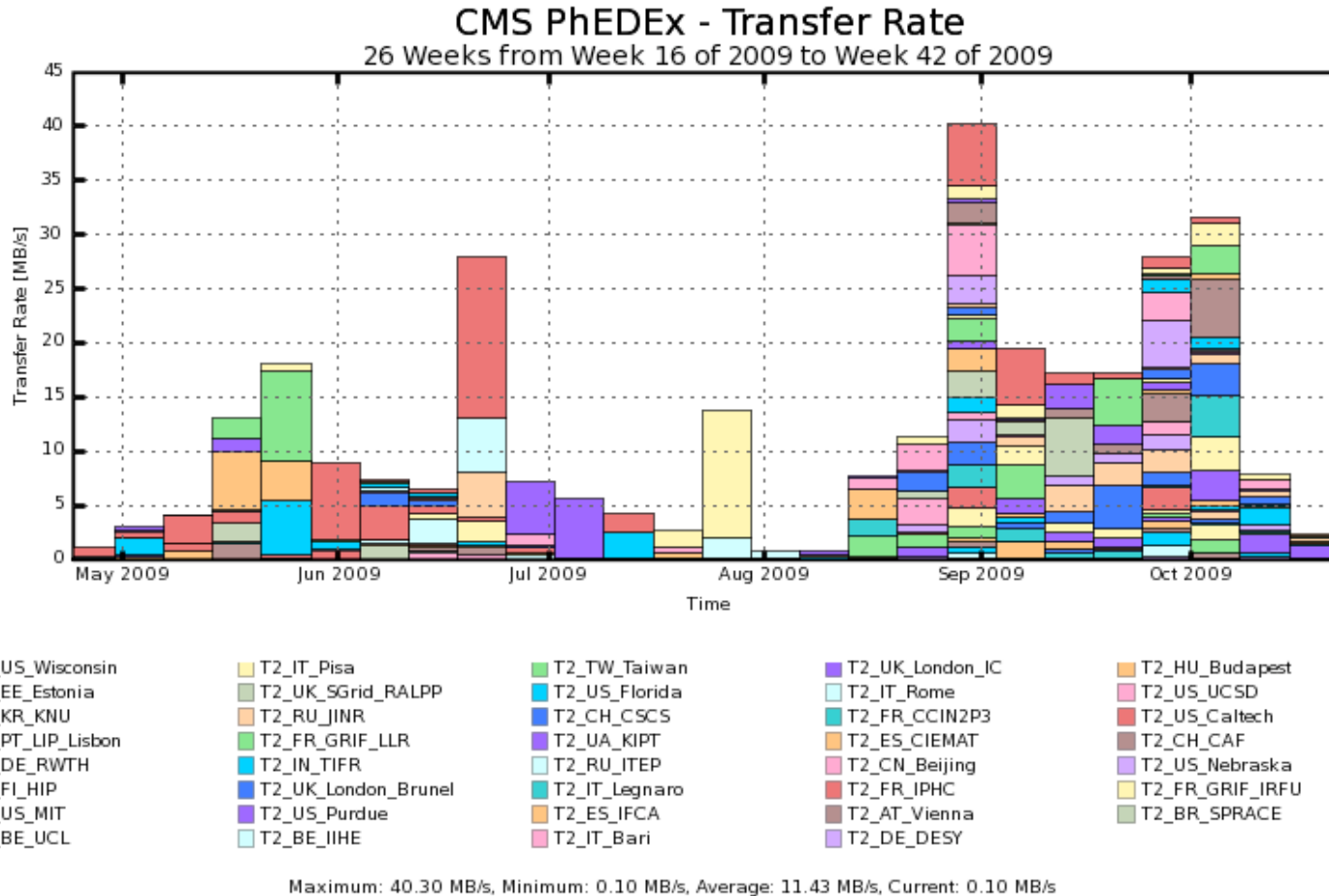Maximum: 211.91 MB/s, Minimum: 0.00 MB/s, Average: 72.21 MB/s, Current: 0.50 MB/s

CMS PhEDEx - Transfer Rate
14 Days from 2009-06-01 to 2009-06-15

**IN2P3 -> T1s** in /Prod

[STEP'09]

T1_FR_CCIN2P3_Buffer to T1_IT_CNAF_Buffer
T1_FR_CCIN2P3_Buffer to T1_ES_PIC_Buffer
T1_FR_CCIN2P3_Buffer to T1_TW_ASGC_Buffer
T1_FR_CCIN2P3_Buffer to T1_UK_RAL_Buffer
T1_FR_CCIN2P3_Buffer to T1_US_FNAL_Buffer
T1_FR_CCIN2P3_Buffer to T1_DE_FZK_Buffer

Maximum: 122.06 MB/s, Minimum: 0.00 MB/s, Average: 50.07 MB/s, Current: 0.39 MB/s

# Transfers: IN2P3 -> T2's



CMS PhEDEx - Transfer Rate
26 Weeks from Week 16 of 2009 to Week 42 of 2009

Maximum: 40.30 MB/s, Minimum: 0.10 MB/s, Average: 11.43 MB/s, Current: 0.10 MB/s
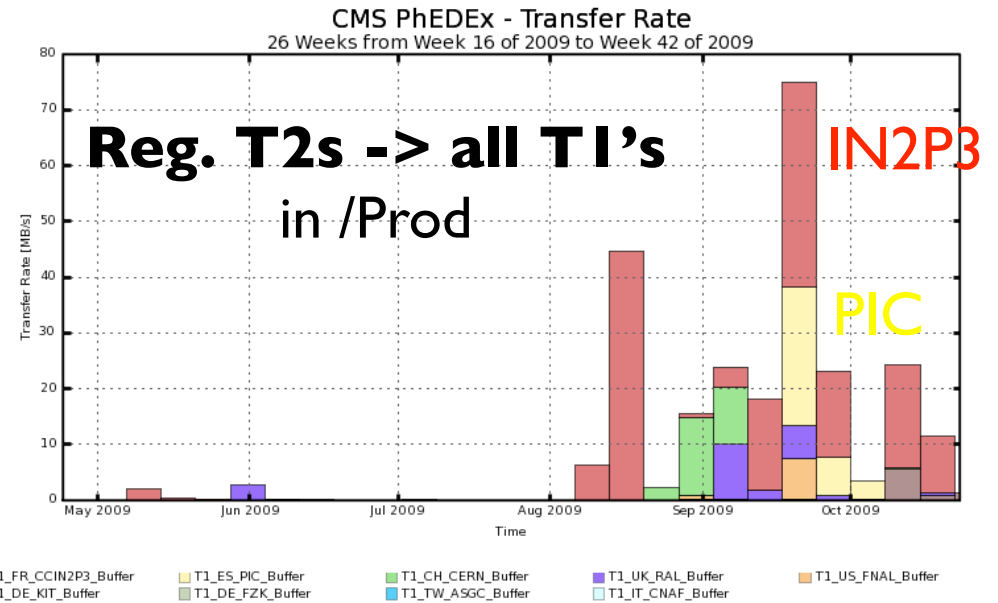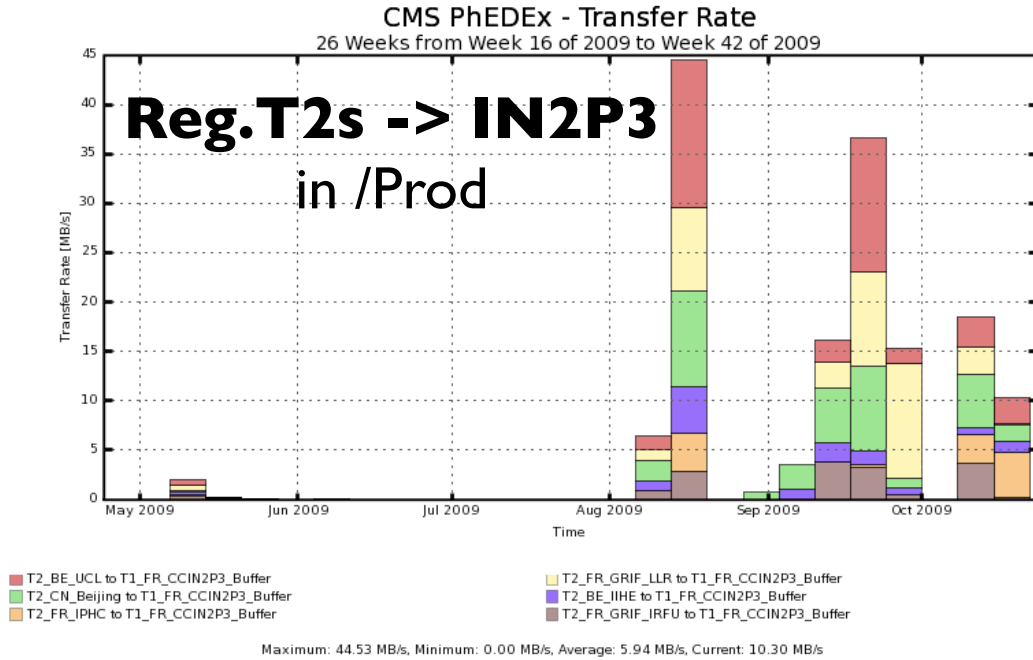
Constant activity since September

Target in CCRC08 was ~80 MB/s averaged over a long period

✦ still below, despite lots of data custodial at T1_FR_CCIN2P3 (~600 TB).

DataOps scheduled a round of 'DDT-style' tests IN2P3->T2_* last week to measure export rates

# Transfers: T2's -> IN2P3 and other T1's



**Reg.T2s -> IN2P3** in /Prod



**Reg. T2s -> all T1's** in /Prod

Assuming CCRC'08 targets:

- ✦ the MC production rate from T2s in the France/Belgium/China region averaged over a long period should be 7.4 MB/s

We are way higher than that after the summer

# Link commissioning status

http://lhcweb.pic.es/cms/CommLinksReports/CommissionedLinks_Sites.html

**T1_FR_IN2P3:**

✦ Export links commissioned, except for some in T2_RU/T2_TR region (rate limitations)

  - http://cmsweb.cern.ch/phedex/prod/Components::Links?from_filter=T1_FR&andor=and&to_filter=.*&Update=Update#

✦ Import links OK, also some non-regional links (not all of them, though)

  - http://cmsweb.cern.ch/phedex/prod/Components::Links?from_filter=&andor=and&to_filter=T1_FR&Update=Update#

**T2's in France/Belgium/China region:**

✦ All fully equipped with downlinks and with many backup uplinks

✦ T2_FR_CCIN2P3 exports still inactive during namespace migration

✦ Remarkably, T2_FR_GRIF_LLR also has lots of T2<->T2 links

# Ops efficiency and Communication

A good coverage of CMS Ops includes:

- ✦ Fulfill your site contact responsibilities
  - Good summary in DataOps slides (next talk)
- ✦ Attend regularly the Ops weekly meetings
  - Provide the brief weekly report every Monday to FacOps
  - Come prepared and discuss current issues on SAM, JR, ... in full depth
  - Give feedback to DataOps on production activities
- ✦ Give complete and precise answers to questions by FacOps and DataOps
  - Meetings, HN, private communications, ...
- ✦ Ask questions yourself !

## Savannah somehow gives a feeling of the rate of issues notifications

- ✦ No Savannah gets opened if a problem is monitored, seen, fixed by CMS contacts onsite <u>before</u> any operator / shifter / user sees it
  - http://snipurl.com/savannah-in2p3
    - IN2P3 102, CNAF 84, ASGC 79, FNAL 74, RAL 48, PIC 30, [ KIT 8 - before: FZK, no history]

## We strongly rely on CMS contacts at T1 sites for efficient operations

# Back-up

# CMSSW deployment

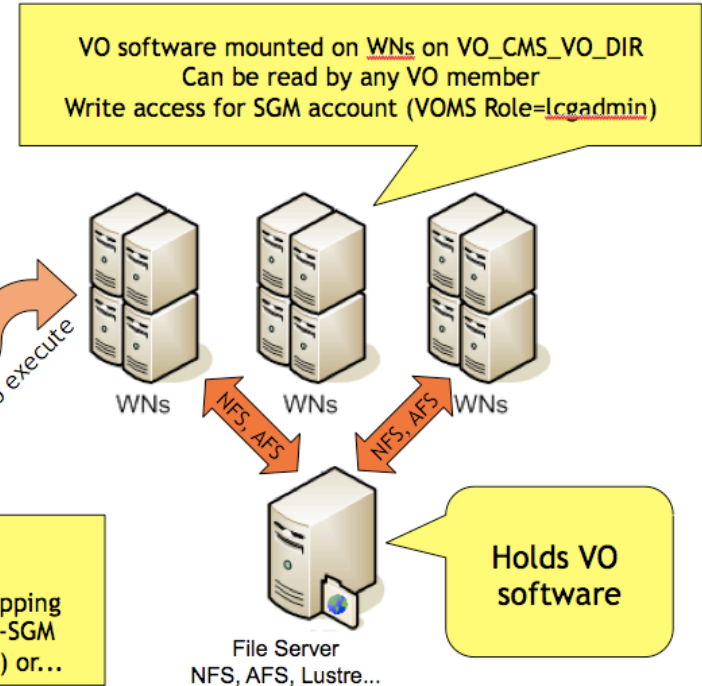## CMSSW installed via Grid job on EGEE and OSG sites

✦ Basic strategy: use RPM (with apt-get) in CMS SW area

### CMSSW_314 deployment

- release announcement on Saturday Oct 3rd at 13h22
  - first installation jobs submitted at 13h41
- status: CMSSW_314 release deployed and ready for Oct-X start-up on day
  - EGEE: submitted to 51 Computing Elements (CE), 44 were DONE after few hrs (see plot); started to follow up on tails over the weekend already
  - OSG: release not tagged into the tag collector, so installed manually; smoothly and quickly completed in most OSG T2/T3

_Credits_:
CMSSW deployment teams in Facilities Ops:
  Bockjoo Kim,
  Joris Maes,
  Lukas Vanelderen,
  Petra Van Mulders,
  Ilaria Villella,
  Christoph Wissing

EGEE

CE Info provider:
Publishes installed software (Tags)

Job execute

VO software mounted on WNs on VO_CMS_VO_DIR
Can be read by any VO member
Write access for SGM account (VOMS Role=lcgadmin)

WNs    WNs    WNs

NFS, AFS    NFS, AFS

Gatekeeper:
- User certificate check
- VOMS role dependent mapping
- cms-USER, cms-PRD, cms-SGM
  or CMS-DE (German sites) or...

CE

Holds VO software

File Server
NFS, AFS, Lustre...

## On EGEE and OSG:

✦ CMSSW releases get routinely installed smoothly in most sites within few hrs from the release

OSG | Year 2009 | |
--- | --- | ---
**Installed Releases** ( 2_2_4 TO 3_3_0 ) | | 33
**Total number of CEs** | | 25
**Total Installations** | | 621
**Total Removals** | | 365

_Credits_: Bockjoo Kim

# **Ticketing systems**

## GGUS

## Savannah



# GGUS

✦ Long tradition of the standard Global Grid User Support system

- Reaches the WLCG site-admins and the fabric-level experts

# Savannah

✦ Problem tracking, troubleshooting reference, statistics, …

- Reaches 'squads' easy to define: CMS contacts at Tiers, tools/services experts, …

- More: baseline tool for Offline Computing shifts, integrated with other CMS projects, …

# Ticketing systems

## Wouldn't a single ticketing system be preferable?

✦ Of course. BUT: is there one with all the features CMS uses for Ops?

## CMS requested a Savannah-to-GGUS bridging

✦ Work finalized. Now ready to be used. Start soon to gain experience in Ops

- Thanks to Guenter Grein (GGUS), Yves Perrin (LCG/SPI) and Simon Metson (CMS) for their great efforts in the technical implementation and testing

# STEP'09 :: IN2P3

## Pre-staging started on June 8-12th due to scheduled HPSS upgrade

- ✦ Site-operated pre-staging approach was chosen *(1)*
- ✦ HPSS v.6.2 interfaced to TReqs interface was used
  - files sorting based on the file position on tape

## Sizable multi-VO activity throughout STEP'09

## High loads observed on HPSS (June 8-13th) *(2)*:

- ✦ Due to all CMS activities simultaneously, in particular CMS analysis at the T2, and also other VOs activities
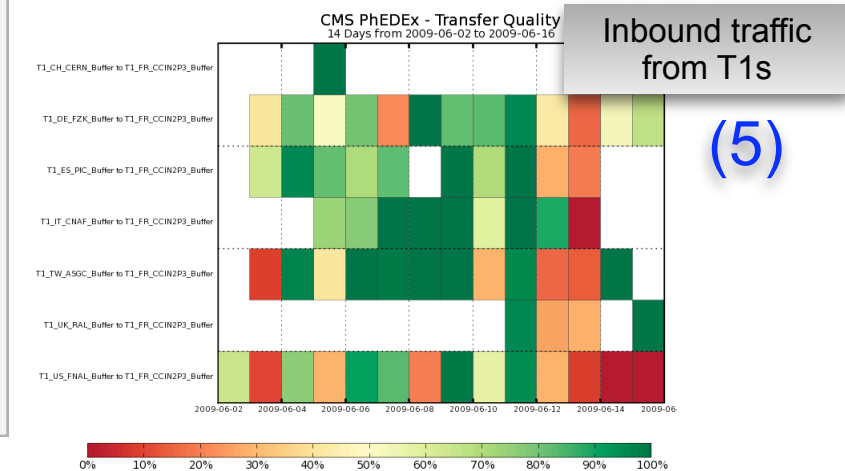  - Decided to suspend T2 analysis activity during STEP'09

## Reprocessing

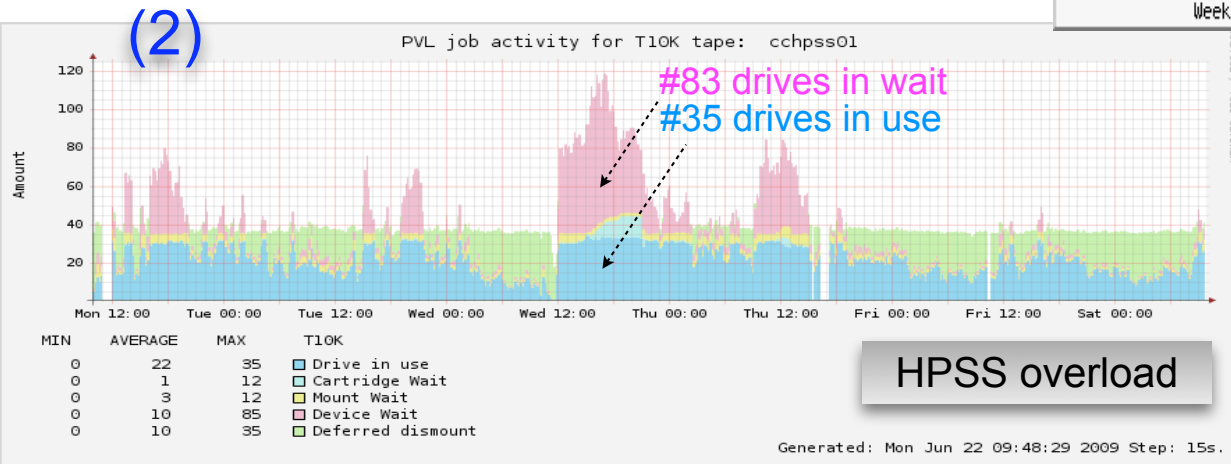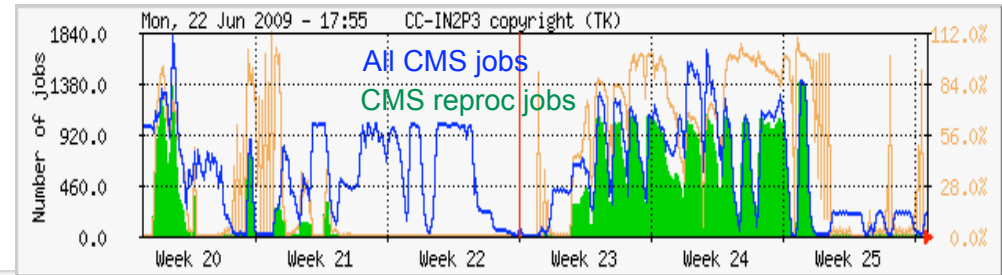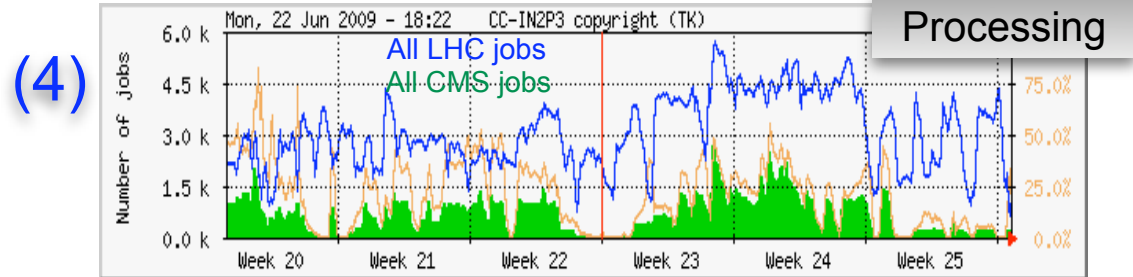- ✦ High reprocessing load by CMS and other VOs *(4)*
  - Failures mainly due to stage-out
- ✦ File distribution per tape on a typical day averages at ~10 *(3)*

## Transfer

- ✦ Relatively smooth
  - some structure (in quality) to be cured, mainly in T1-T1 *(5)*

# STEP'09 :: IN2P3 in plots



(1) Tape: bySC-raw:cms — HPSS->dCache migration

(2) PVL job activity for T10K tape: cchpss01 — HPSS overload
#83 drives in wait
#35 drives in use

(3) File distrib on tapes — nombre de fichiers par cartouche cms-day04-prestage

(4) Processing — All LHC jobs / All CMS jobs / CMS reproc jobs

(5) Inbound traffic from T1s — CMS PhEDEx - Transfer Quality

T1_FR_CCIN2P3

# STEP'09 :: IN2P3 *[pre-staging]*

| | Site-operated | Central SRM script | PhEDEx agent |
|---|---|---|---|
| ASGC | | | X |
| CNAF | | X | *Castor + StoRM: need additional work in PhEDEx* |
| FNAL | X | | |
| FZK | X | *Tape issues: preferred manual* | *Tape issues: preferred manual* |
| IN2P3 | X | *HPSS downtime on week-1: preferred manual* | *HPSS downtime on week-1: preferred manual* |
| PIC | | | X |
| RAL | | | X |



~105 MB/s

Scheduled HPSS down

52 MB/s **Target**

| | Target [MB/s] | 2-Jun | 3-Jun | 4-Jun | 5-Jun | 6-Jun | 7-Jun | 8-Jun | 9-Jun | 10-Jun | 11-Jun |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ASGC | 73 | *Digesting migration* | 140 | 170 | 190 | 160 | 145 | 150 | 140 | 150 | 220 |
| CNAF | 56 | 380 | 300 | 160 | 240 | 240 | 270 | 105 | 80 | 125 | 240 |
| FNAL | 242 | 280 | 200 | 200 | 120 | *Still staging previous day* | *Recovering from backlog* | | 379 | 380 | 400 |
| FZK | 85 | *Tape system not available [unscheduled downtime]* | | | | | | | *Participated in pre-staging but performance not clear* | | |
| IN2P3 | 52 | *Tape system not available [scheduled downtime]* | | | | | | | 96 | 99 | 120 | 103 |
| PIC | 50 | 60 | 61 | 106 | 83 | *Samples not purged* | *Samples partially on* | 99 | 142 | 123 | 142 |
| RAL | 40 | 250 | 230 | 160 | 140 | 135 | 190 | 170 | 100 | 220 | 180 |

# STEP'09 :: IN2P3  *[CPU efficiency]*



An example day: **June 11th**
[daily plots collected throughout STEP'09]

IN2P3
[ dCache HPSS ]

**with pre-staging**

**without pre-staging**

Previously failed jobs
might have already
triggered the pre-staging

# Measured every day, at each T1 site. Mixed results:

✦ Very good CPU efficiency for **FNAL**, **IN2P3**, (**PIC**), **RAL**

✦ Not so good CPU efficiency for **ASGC**, **CNAF**

✦ Test not significant for **FZK**

# Current understanding:

✦ Test demonstrated the significant effect of pre-staged data for processing

✦ Site specifics to be investigated: **IN2P3 not one of these**