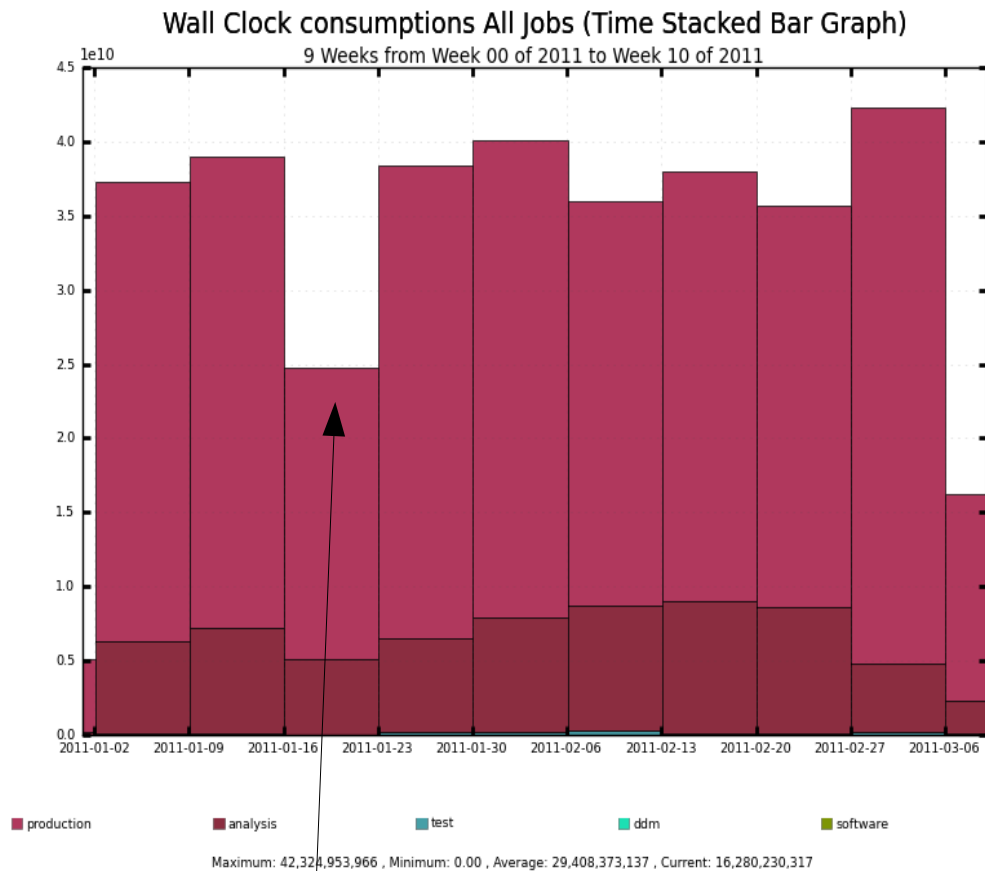# ATLAS Operation report

# GDB

**S. Jézéquel**

9 March 2011

- **Operation report (Dec 2010- Mar 2011)**

- **Short term implementation : Gradual breaking of cloud model**

- **ATLAS Data Distribution model : 2011**

# Operation activities(Dec 2010-Mar 2011)

- No specific activity during this period

- Consolidation (ATLAS+sites) and preparation for 2011 activities :
  - Service certificate replace Kors certificate for DDM activities
  - Oracle upgrade at CERN + ATLAS database splitting
  - Validation of EOS technology at CERN
  - Taiwan DISK : Castor → DPM (to be finished)
  - Setup automatic export of ID calibration data from T0 to TAIWAN/Valencia
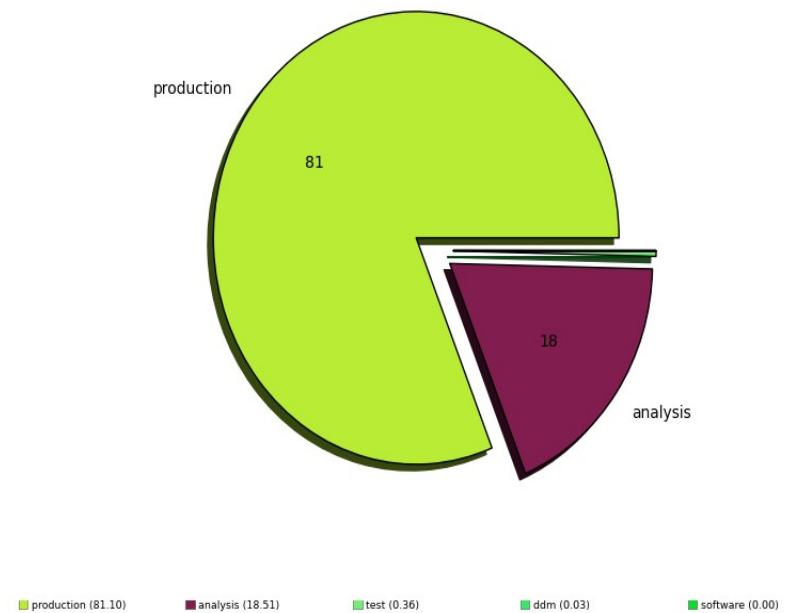
- 2010 Heavy-Ion reprocessing starting this week

9 March 2011

# Processing activity: Jan-March 2011

## All ATLAS Grid activity (January-March 2011)



Wall Clock consumptions All Jobs (Time Stacked Bar Graph)
9 Weeks from Week 00 of 2011 to Week 10 of 2011

Maximum: 42,324,953,966 , Minimum: 0.00 , Average: 29,408,373,137 , Current: 16,280,230,317

production    analysis    test    ddm    software

~66 k full CPUs

Wall Clock consumptions All Jobs (Pie Chart in percentage) (Sum: 100.00)

production 81

analysis 18

production (81.10)    analysis (18.51)    test (0.36)    ddm (0.03)    software (0.00)

CERN Oracle upgrade : 16 Jan (LCGR) -17 Jan (ADCR+ATLR)

4

## Production
## (MC + group production)

## Analysis
## (user+phys. groups)

Wall Clock consumptions All Jobs (Pie Chart in percentage) (Sum: 100.00)



Wall Clock consumptions Good Jobs (Pie Chart in percentage) (Sum: 100.00)



- US (19.93)
- DE (17.21)
- FR (12.98)
- UK (11.41)
- ND (6.92)
- NL (6.29)
- CA (5.77)
- CERN (5.76)
- ES (5.50)
- IT (5.09)
- TW (3.15)

- US (34.60)
- DE (16.09)
- FR (12.57)
- CERN (8.08)
- IT (7.73)
- UK (7.48)
- NL (3.92)
- CA (3.90)
- ND (2.30)
- ES (1.95)
- TW (1.38)

9 March 2011

# Storage news

### Merging of space tokens DATADISK/MCDISK

- Data migrated with DDM (FTS transfer + central deletion)
- Only primary replicas were transfered (minimize activity)
- Exception : RAL : Internal Castor migration
- Timescale
  - January : T2/T3 (no other major transfer activity)
  - February-March : T1 (no other major transfer activity)
  - CERN : Will be done when migration Castor → EOS

  Site responsability : When migration is finished, clean remaining dark data

### Space token shares in 2011

- Reference : https://twiki.cern.ch/twiki/bin/view/Atlas/StorageSetUp
- 2011 shares similar to 2010 (validated on 7 March 2011)
- Main changes :
  - Shares ATLASMCDISK and ATLASDATADISK merged
    - → 1 space token host 75-80 % of storage ressources
  - ATLASPRODDISK : 25 TB in T1s (MC production at T1 + ATLAS managed stagin buffer)

# Operation issues

- **RAL :**
  - **Problem with Castor during migration from bad servers (Christmas period)**
  - **MCDISK was full because of bug in ATLAS central deletion**
    - → **No MC production in UK cloud during a week**

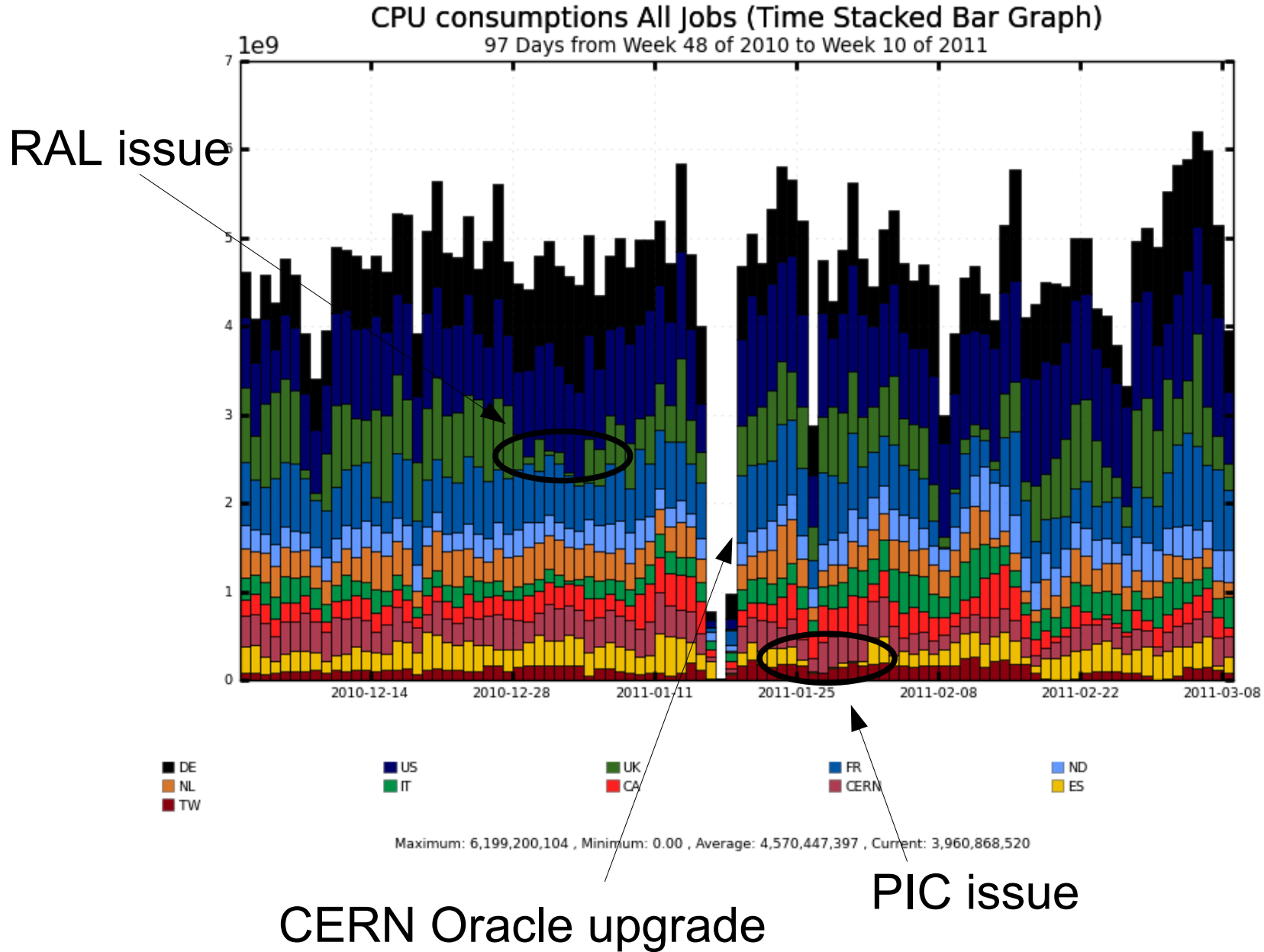- **PIC : Temporary lost 800k files/250 TB after file system corruption (GGUS : 66409)**
  - **Took a week to get list of unaccessible files**
  - **PIC announced that most of files could be accessible after internal migration**
  - **Too many files to recover from outside (1/3 was available outside)**
    - → **Data recovered/consolidated within PIC (~1 week)**
    - → **No production in ES cloud during 2 weeks + some unaccessible data**

- **Castor/srm upgrade to 2.10 (CERN) (GGUS :):**
  - **Was not able to rollback after experiencing instabilities**
  - **Solved within a day by patching 2.10**

- **Thousands of jobs from same user using > 4 GB memory : Solved by contacting user**

RAL issue

CERN Oracle upgrade

PIC issue

8

# Breaking cloud model

- ATLAS wants to break the cloud model to get more flexibility

- Obvious constraint : Should match the network connectivity between sites

- Done step-by-step :
  - To ensure that scalibility issues can now be overcome
  - Adapt the monitoring tools
  - Train ATLAS shifters : who is responsible in case of problem

- Current actions :
  - Prepare LFC consolidation at CERN
  - Some T2s running G4 simulation with input/output files transfered from/to different T1s
  - Direct transfers between some T2s and all T1s

- Future actions :
  - Promote 'good' T2s to host primary replicas (only in T1s today)

9 March 2011

# LFC consolidation at CERN

- **Goal :**
  - All LFCs agregated in a single LFC at CERN
  - Read-only replica in another site (probably BNL)

- **Reason :**
  - ATLAS experienced LFC downtime over few weeks (Summer 2010)
  - Current LFC model: single point of failure
    - $\rightarrow$ all stored data within cloud can be unaccessible

- **Status :**
  - Discussion between ATLAS and WLCG/CERN to validate the merging procedure
  - Identifiy possible inconsistencies between catalogs before merging

- **Timescale :**
  - One LFC migration at a time (should be done within days)
  - Expected to be done during spring/summer 2011

9 March 2011

# Cross-cloud production

◆ **Reason :**

- ◆ **Allocate more CPU ressources for urgent/big simulation tasks**

  → **Gather CPUs from sites outside the cloud**

- ◆ **Avoid inbalance between T1 storage ressource and CPU ressource within cloud**

- ◆ **Continue to use T2s CPUs when T1 SE is down**

◆ **Necessary connectivity :G4 simulation jobs in T2s with output transfered to T1**

**1000 cores : Transfer rate ~ O(1) MB/s : Easy**

◆ **Main issue :**

- ◆ **Setup FTS channels to avoid to go through STAR-T1/T2 channels (T2D topic)**

- ◆ **Adapt monitoring**

◆ **Current situation**

- ◆ **Some big T2 sites already associated to many Tier1s**

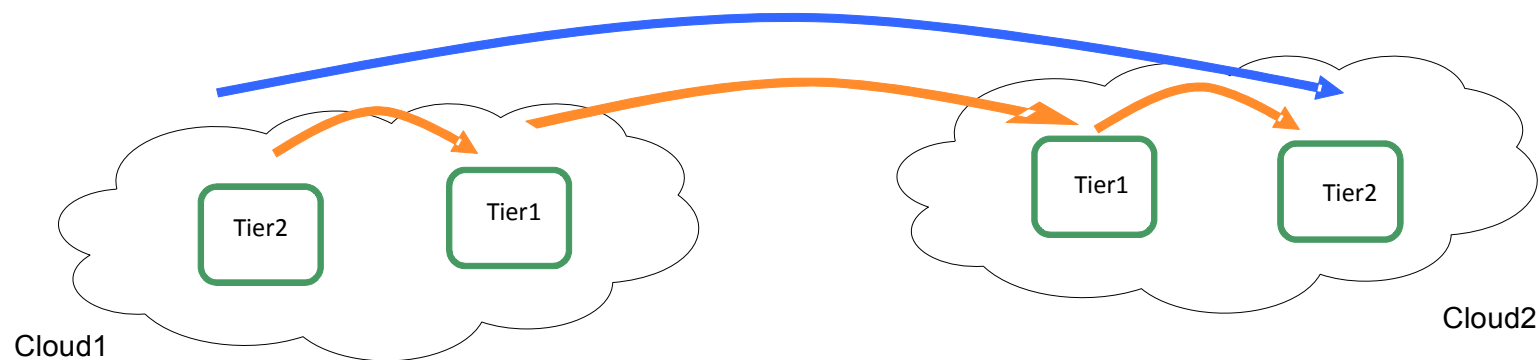| NL sites | Pilots | Latest | defined | assigned | waiting | activated | sent | starting | running | holding | transferring | finished | failed | cancelled |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| NL | 2207 | 02-14 21:40 | 0 | 260 | 144 | 4207 | 0 | 11 | 4057 | 73 | 3162 | 831 | 241 | |
| ALL | | | 0 | 260 | 144 | 4207 | 0 | 11 | 4057 | 73 | 3162 | 831 | 241 | 70 |
| DESY-HH | 445 | 02-14 21:40 | 0 | 0 | 0 | 273 | 0 | 1 | 175 | 0 | 817 | 1 | 14 | 0 |

11

# Data collection into T2s

**Current DDM model**



ISSUES T2→ T1 for T2S WITH BAD CONNECTIVITY

**DDM model : Version +1 (Under validation in Italy cloud)**



2 possible paths → Need monitoring to optimise path (Depend on file size)

Expected to use direct transfers for small files

9 March 2011

# T2 connectivity : Monitoring

◆ **Extend current channel validation (T0→T1, T1a→T1b, T1a→T2a)**

**to T1a→T2b and T2a→T2b**

**http://bourricot.cern.ch/dq2/ftsmon/sonar_view/cached/**
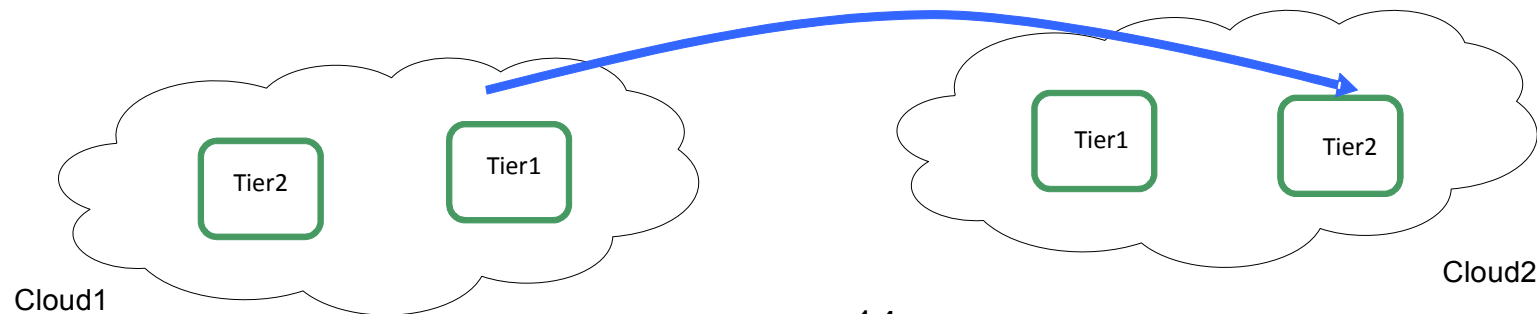
# Data collection into T2s (2)

- **Issue :**
    - **T2s with big storage capacity collect lots of data**
    - **Current model : cross-cloud transfers go through T1s**
        - **Reason : Good connectivity between clouds through T1s/LHCOPN**
    - **But :**
        - **Triggers additional activity in T1s SE (copy+delete)**
        - **Transfers to T2s are stuck if T1 SE is down/full**
        - **Transfers can be delayed if huge activity for T1 transfer**
- **Target :**
    - **Minimize useless load on T1 SE**
    - **Minimize intermediate steps → Less sensitive to intermediate site availability**
    - **First use case : Collect group production at T1s into group storage at T2**

**DDM model : Version +1' (Implemented)**



Cloud1  Tier2  Tier1  Tier1  Tier2  Cloud2

9 March 2011

# Direct cross-cloud T2 connectivity

- Select good T2 sites which will always transfer from/to all T1s
  - T2Ds
- T2D current list :
  - All US T2s
  - DESY-HH,DESY-ZN, GRIF-LAL,GRIF-LPNHE, INFN-NAPOLI, IFIC
- A long list of sites under probation
- ATLAS would like to have as many sites as possible
- Triggered many network studies (UK for example)
  - → More sites will be added soon

  **LHCONE will be a key component for this policy**
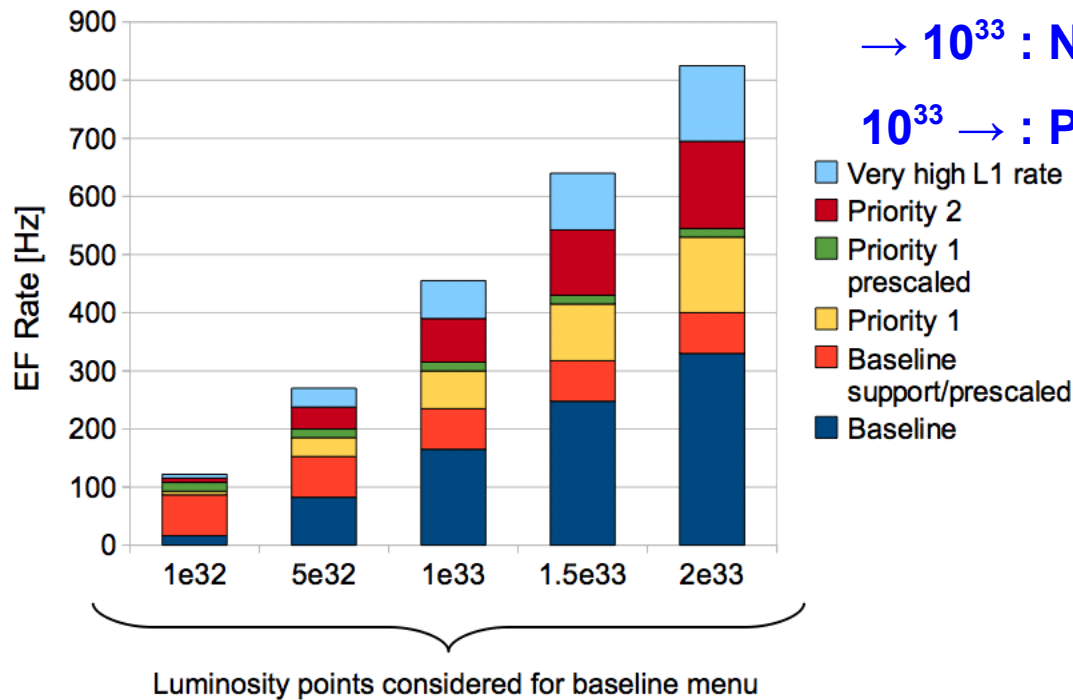
- Request additional FTS channels :
  - T1s → T2D
  - T2Ds → T1 (currently go through STAR-T1 channel)
  - Implementation/optimisation discussed in T1 Coordination meeting
  - Triggered discussions with FTS dev.

9 March 2011

# T2 connectivity : Summary

+ **Full multi-hop model under validation**

  + **Hopefully generalised in March 2011**

+ **Monitoring cross-cloud transfers : Done**

  + **Migration to well supported monitoring framework under way**

+ **DDM connectsT2Ds to all T1s : Done**

  + **First list of T2Ds defined**

  + **FTS channel setup being optimised**

9 March 2011

# 2011 Data Distribution Model

New requirements from ATLAS physic/trigger community to collect more data

→ Increase mean Event Filter mean rate to 400 Hz (limitation is T0 capacity)

→ $10^{33}$ : No additional cut

$10^{33}$ → : Priorities defined to limit to 400 Hz



Legend:
- Very high L1 rate
- Priority 2
- Priority 1 prescaled
- Priority 1
- Baseline support/prescaled
- Baseline

**Computing issues**

Transfer rate of fresh data from CERN

Storage capacity for primary replicas produced over year

→ Review of the ATLAS Data Distribution Model

17

9 March 2011

# 2011 Data Distribution Model

- Reduce RAW size : zip files

  - Gain a factor 2 (many empty calorimeter cells)

    - → Compression factor close to 1 when writing on TAPE

  - Zipping is done at T0 level

  - Unzipping files is done during file reading (< 0.1 s  in addition)

  - Will be implemented in coming days

- If needed, transfer CERN CAF CPU resources to T0 from prompt data reco.

  - → CERN cannot replace a stuck T1 for reprocessing campain

  - → All Tier1s should reach promised availability

- In addition to TAPE copies,

  - 1 RAW copy on DISK to allow prompt access for 'discovery' studies

    (few to few 100k events accessible within 24 hours)

9 March 2011

# 2011 Data Distribution Model (2)

- **'Life without ESD' :**

  - **Restrict number of ESD replicas (2 copies)**

    → **Promote analysis from AOD/DESD**

  - **Lifetime of 6-8 weeks for bulk ESD streams**

- **AOD/DESD:**

  - **Number of replicas adapted to available DISK ressources**

  - **Promote big T2Ds to host primary replicas (+10 TB)**
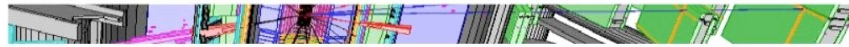
    → **Include them in schedulded downtime coordination ?**

**LHC one : key component to deliver datasets from these T2s**

9 March 2011

# 2011 Data Distribution Model: Rate

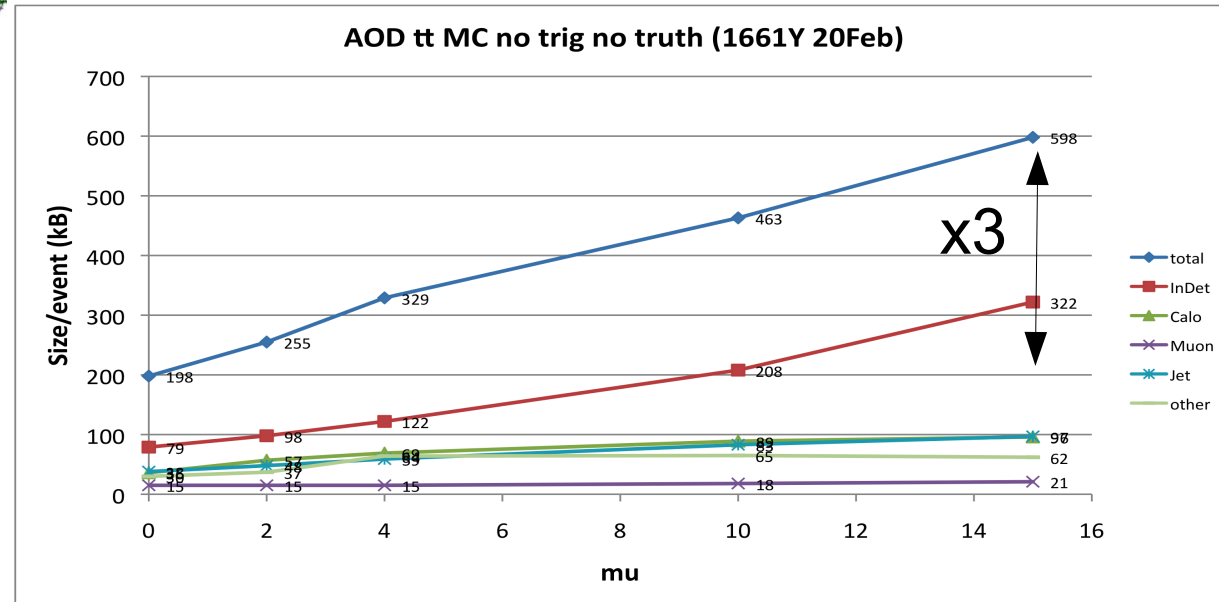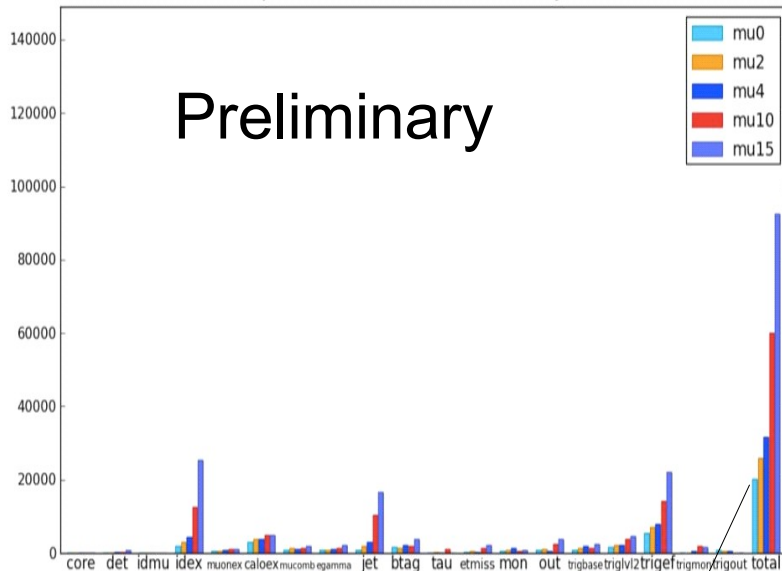- Transfer rate T0 → T1s for 400 Hz

  - 2 RAW copies : 2*320 MB/s (0.8 MB/evt)

  - 1 ESD copy : 640 MB/s (1.6 MB/evt)

  - 2 AOD +2 DESD copies : 4 * 110 MB/s (0.275 MB/evt)

    - → Total : 1720 MB/s (does not include internal read/write within T0)

- Transfer activity will be smoother : No run validation before exporting data

9 March 2011

➢**LHC luminosity increase ← including more protons per bunch**

→    **Number pile-up events increases rapidly with luminosity**

➢**Pile-up beginning/end 2011 (L ~O($10^{32}$)/O($10^{33}$))  ~0/15**



**Less CPU ressources for CERN CAF activities → CERN not backup for T1s reprocessing**

•Pressure on software developpers to reduce the size increase → Improvments each day

Another potential source of tension for computing ressources
      → Monitoring and reactivity will be necessary

# Conclusion

- **Smooth ATLAS computing activity during the last 3 months**
  - **Includes usual rate of site/tools issues**

- **Next months will be much more challenging than 2010**
  - **All consolidation activities should proceed quickly**
  - **Sites (especillay T1s) should keep up with expected availability**

- **Many medium/long term developments also coming hopefully**
**in collaboration with WLCG and other LHC experiments**

9 March 2011