

## **T2s/T3s consolidation**

**S. Jézéquel (LAPP)**

**ATLAS Computing Jamboree  
5 March 2018**

- T3s and diskless T2s policy :

- Discussed and endorsed in 2017 by International Computing Board (ICB)

→ You should be aware of the points presented today

- Goal : Optimise the manpower (ADC/site admin/squad) and computing activities

vs

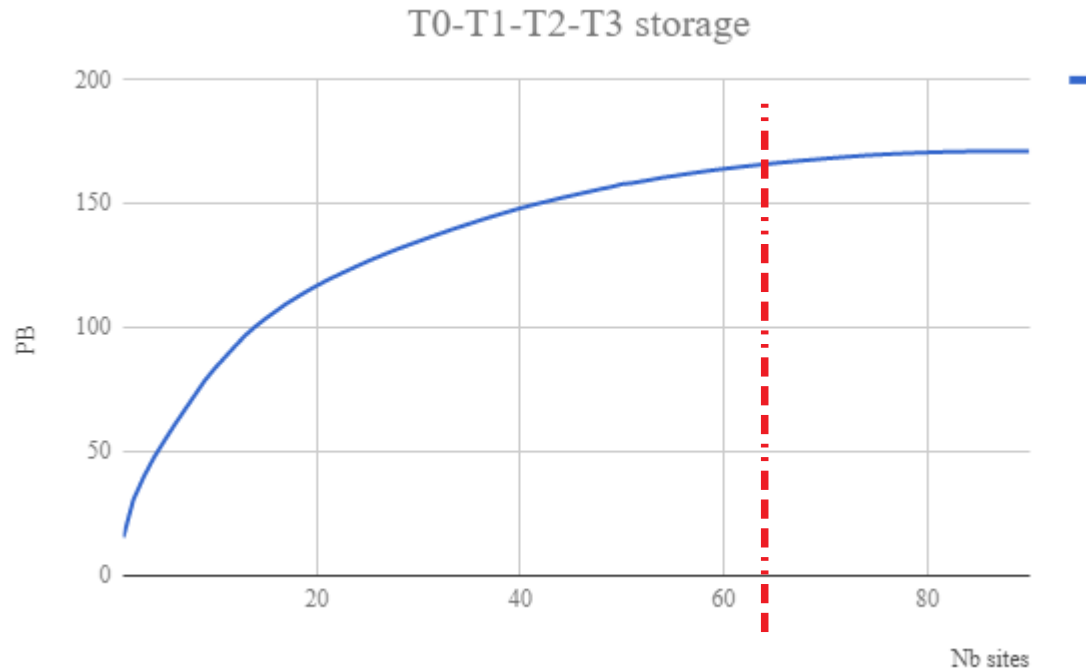
amount of computing resources

- More general presentations tomorrow (Xin,...)

~2013 : Proposition by E. Lancon as ATLAS Computing coordinator

**'ATLAS recommends sites with small storage  
to focus hardware investment in CPUs instead of storage'**

- Storage maintenance is dominated by Grid layer upgrade/maintenance instead of scalability issues
  - Support from Rucio support per site is roughly independant from storage size
  - Observation : Most of unreliable T2 sites are among the small ones
  - 10 % of T2 storage (8 out of 80 PB) represents 50 % of DATA/SCRATCH Rucio endpoints
- → Set limit on recommendation to 0.5 % of T2 storage : 400 TB in 2017
- **T2 consolidation or T2 Diskless policy**



400 TB  
limit

- Should not be considered as a call for one site per country or FA
  - Grid infrastructures allowed to get more computing resources than restricted to national computing facilities
  - But provide reliable sites (hardware and manpower)

Maximise computing resource delivered to ATLAS Collaboration  
at a reasonable price in manpower/hardware cost

Example : T2s in Canada

- Implementation :
  - Jobs in Site A read/write files from site B storage
  - Configuration defined in AGIS
- Early 2017 : Validation of setup and operation stability (E. Vamvakopoulos IN2P3-CC)
  - 2 romanian sites (800 cores total) were configured to access storage in 3rd romanian site ( 1.2 PB)
  - Motivation : Difficulty to maintain stable storage over long time  
→ diskless candidates
  - Performance after diskless implementation :
    - ◆ No significant degradation of site production efficiency on both sides

- Summer/fall 2017 :
  - Which sites :
    - Storage < 100 TB and in second step <400 TB
    - Network connection to a possible partner > 1 Gb/s
  - Date of implementation :
    - Immediately : remove problematic storages
    - In coming months/years (match country policy on decommissioning)
  - Associated ICB rep. were asked to define timeline :
    - Most aware of funding/manpower optimisation
    - With sites and squads
  - T2 Diskless policy
    - presented in ADC and ICB meetings
    - Endorsed by ICB Representatives in fall 2017

- Immediately implemented on 6 sites
  - Romania (2), Sweden (1), Russia (3), Austria (1)
- Few countries will implement it inline with their decommissioning policy (few storages in 2018)
- In some country, keeping storage at each site is strategic
- No solution evaluated for isolated sites

| RSE  | Total<br>(storage) ▼ | Used (rucio) | Used (other) | Quota (other) | Used            |
|--|----------------------|--------------|--------------|---------------|-----------------|
| <a href="#"><u>LUCILLE_DATADISK</u></a>                | 84.32                | 30.07        | 0.00         | 0.00          |                 |
| <a href="#"><u>TR-10-ULAKBIM_DATADISK</u></a>          | 131.94               | 113.84       | 0.00         | 0.00          |                 |
| <a href="#"><u>CA-SCINET-T2_DATADISK</u></a>           | 140.00               | 120.00       | 0.00         | 0.00          |                 |
| <a href="#"><u>NCG-INGRID-PT_DATADISK</u></a>          | 177.00               | 150.34       | 0.00         | 0.00          |                 |
| <a href="#"><u>EELA-UTFSM_DATADISK</u></a>             | 180.00               | 161.22       | 0.00         | 0.00          | T2s             |
| <a href="#"><u>UKI-SOUTHGRID-CAM-HEP_DATADISK</u></a>  | 207.00               | 188.07       | 0.00         | 0.00          |                 |
| <a href="#"><u>UTA_SWT2_DATADISK</u></a>               | 225.00               | 193.87       | 0.00         | 0.00          | March 18        |
| <a href="#"><u>UKI-SOUTHGRID-BHAM-HEP_DATADISK</u></a> | 225.40               | 205.03       | 0.00         | 0.00          |                 |
| <a href="#"><u>BEIJING-LCG2_DATADISK</u></a>           | 310.00               | 281.95       | 0.00         | 0.00          | Total (~3.5 PB) |
| <a href="#"><u>RRC-KI_DATADISK</u></a>                 | 315.56               | 243.71       | 0.00         | 0.00          |                 |
| <a href="#"><u>PSNC_DATADISK</u></a>                   | 340.85               | 305.78       | 0.00         | 0.00          |                 |
| <a href="#"><u>RO-02-NIPNE_DATADISK</u></a>            | 346.35               | 313.20       | 0.00         | 0.00          |                 |
| <a href="#"><u>WEIZMANN-LCG2_DATADISK</u></a>          | 360.00               | 316.71       | 0.00         | 0.00          |                 |



- Implementation be pursued over coming years
- Keeping 0.5 % limit → Threshold +15 % each year (flat budget + cost reduction)
  - 460 TB in 2018, 600 TB in 2020, ....
- First step towards data lake @ HL-LHC ?

- Aim : Optimise the ADC support devoted to T3s
  - Most of the unreliable ATLAS sites are T3s
    - Few site admins and possibly fast turnaround (→ training issue)
    - Some T3s are built with recycled hardware
      - In 2017 : 6 % of ADC central operation for 0.5 % of computing resources


Recommendations not meant for :

- HPC/Cloud sites : Still require dedicated manpower but worth doing it since big
- T3 colocated with T1/T2 (benefit from local Grid expertise)
- SLAC : T2 untill recently → Still large infrastructure and team
- In general reliable T3s with a strong Grid support team

● Recommendations should be applied on :

- new sites (~ 1 request per month)
- sites changing their configuration

- T3 wishes :
  - Contribution to ATLAS Grid through CPUs
    - **Site offers resources to ATLAS**
    - Boinc, US-Connect sites, Greece, South Africa, Geneva,...
    - In the past : Argentina, Armenia
  - Installing LOCALGROUPDISK :
    - **Site benefits from ATLAS Grid infrastructure to collect files from Grid**
    - Few US sites, Greece, South Africa, Turkey
    - In the past : US, Armenia, Argentina, Switzerland

- No pledge resources and no reliability commitment
  
  - ADC priority support : T1/T2 > big T3s > small T3s
    - T3 setup should minimize the source of technical problems and support from ADC experts
    - No critical activity for ATLAS Collaboration affected to T3s
- 
- No primary copy of ATLAS Grid files → No DATADISK
  - No high priority tasks (reliable sites could host such tasks)
  - No Grid analysis queues
- 
- Setting up T2 like (= deploy all services)
    - only if the infrastructure and manpower are already available

- Site with existing Grid infrastructure (CE, cvmfs, squid,...) :
    - Setup a Panda queue as any Grid site (option Event Service)
  
  - Site w/o Grid infra : Target is to minimise complexity compared to a local batch cluster+ storage
    - Avoid Computing Element
    - A priori no Grid storage locally
- Recommend to install ATLAS@home/Boinc or any other lightweight system
- Devt for 2018 : Make visible this contribution at site level for ATLAS accounting

- Only for ATLAS institutes
- Most time consuming point for ADC central:
  - Non-grid storage technology is not widespread among local admins
    - Maintenance relies on very few people and support might be asked to ADC
  - Usually not possible to just clean the storage/catalog
    - Migration campaigns with heavy manual intervention to deal with inaccessible files
- Recommendation for new sites with only LOCALGROUPDISK:
  - Storage > 10 TB (otherwise just use 'rucio get')
  - Install GridFTP door
- Request for all T3 sites :
  - Local team (admin or physicists) responsible to monitor/report issues (ex : Full LGD)
  - New : Prepare data migration/cleaning well before site/SE decommissioning

- ICB rep / squads : Responsibles to check the requirement implementations
  - Faster turnaround of site admin → ICB/squad rôle more important
- To setup the site
  - Setups technical collaboration with some other site admins (at least national coordination)
  - Identified person (ICB rep or squad) should collect the site objectives and adapt them to the manpower available to build and operate the local infrastructure
  - Deployment itself should be tracked in JIRA ticket
- During site operation
  - Publish downtimes (as any Grid site)
  - Respond to GGUS tickets
  - Solve issues within days (storage) or weeks (CPUs) and publish downtime if necessary
- Site decommissioning
  - Identified person (ICB rep or squad) should organise decommissioning few months in advance

- Site exposed in ADC monitoring should have good reliability
  - Short term issues are detected and hidden through HC or FT
  - Broken sites over many weeks with no downtime → Are we missing resources ?
- New ADC policy :
  - Broken site over > 1 month with no concrete action is permanently blacklisted in AGIS
    - Better if handled by squad than ADC
  - Each year, permanently blacklisted sites are reviewed and most probably completely decommissioned in AGIS and Rucio (ICB rep + AGIS site contact informed)



- Twiki (link) : the recommendations to decide the optimal T3 configuration and operation
- Tried to apply rules on obvious issues
- Removal from AGIS : T3 broken sites since many years : 4 sites (AM-04, ROMA2, GR-01, UJ)
- Decomission broken storages since many months : 3 sites (UNIGE, NYU, US2)
- Decomission SCRATCHDISK when no analysis queue : 2 sites (ZA-UJ, GR-12)
- Would be a sign of success : Shifters would accept to report T3 storage issues

- T3s are opportunities for :
  - ATLAS Collaboration to benefit from local CPU resources for production
  - ATLAS users to easily collect data
- But the necessary manpower should be much lower than to operate/support T2/T1 sites
  - 2017 : Implementation of new T3 policy
  - Should continue in 2018 and beyond
  - New sites should be setup accordingly

ICB/squad rôle critical to ensure to that only good sites are integrated in ATLAS Grid infra

# Backup

- RO-16-UAIC : 110 TB, RO-14-ITIM : 30 TB → Romania : 140 TB
  - HEPHY-UIBK (Austria) : 90 TB
  - ITEP : 10 TB, RU-PNPI : 180 TB, RU-MOSCOW-FIAN-LCG2 : 10 TB → Russia : 200 TB
  - SE-SNIC-T2 (Sweden) : 30 TB
- Total : ~ 0.5 PB