



Documenting BW requirements for use cases from the Computing Model at CMS-Tier1s

M. C. Sawley

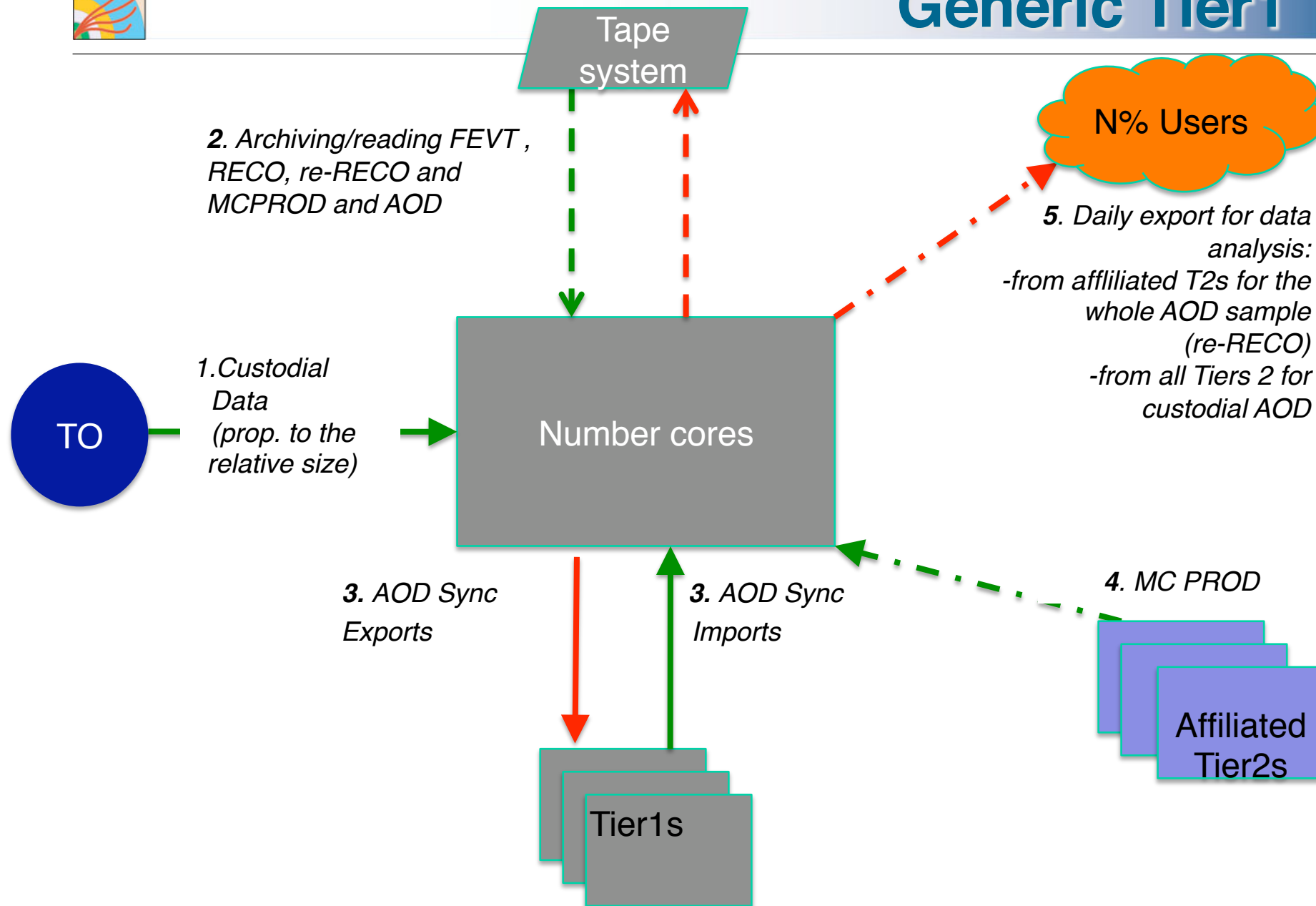
ETH Zurich

12 May 09

Disclaimer: the numbers hereby given are not to be considered final before approval by the CMS Computing Resource Board.



Generic Tier1





Parameters 2009-2010

- **The proposed model evaluates the requirements for data transfer during the 1st data taking period (2009-2010) of 200 days.**
- **100 days of data taking at 20% duty cycle followed by 100 days at 50%**
- **Resources for each site are those pledged for 2009**
- **Raw numbers calculated with no security factor**
 1. **Bandwidth between T0 and T1s during data taking**
 2. **Tape I/O at T1s when writing FEVT, SimFEVT and RECO, and reading RECO, simRECO and re-RECO (in progress)**
 3. **Between T1s for AOD synchronization**
 4. **T2s uploading MCPROD at parent T1**
 5. **T1s exporting selected data for analysis to users**



STARTING POINT

CONSTANTS

(Identical to those taken the CMS resources request for 2009/2010, cf. M. Kasemann))

Trigger rate	300 Hz
RAW Size	1.5MB
SimRAW	2MB
RECO size	0.5MB
AOD size	0.1MB
Total number of events	2590 MEvents
Total number of simulated events	2030 MEvents
Overlap between PD	40%
Total size of RAW	3810 TB
Total size RECO	1270 TB
Total Primary AOD	254 TB



1. Exporting custodial data: methodology

- **T0 → T1s : exporting FEVT**
 - **$BW = (RAW + RECO) \times \text{Trigger frequency} \times (1 + \text{overlap factor})$** . For the chosen parameters, this yields:
 $BW = 2\text{MB} \times 300\text{Hz} \times 1.4 = 840\text{ MB/sec}$, or **6.7 Gb/sec.**
- **Each T1 receives a share according to its relative size in CPUs**
- **Proportional to the trigger rate, event size and Tier-1 size**



2. Tape I/O: methodology

- The maximum rate is computed for concurrent tasks:
 - Writing FEVT
 - Writing MCPROD
 - Writing re-RECO
- Period considered for re-reconstruction: **1 month** (stringent!)
- Proportional to data taking period and inversely proportional to period imposed for re-reconstruction; correlated with MCPROD
- Data rate for writing:
 - $(\text{FEVT (1 month)} + \text{MCPROD (1 month)} + \text{re-Reco+AOD}) / \mathbf{1 \text{ month}}$
- Data rate for reading:
 - $(\text{RECO} + \text{re-RECO} + \text{simRECO} + \text{AOD}) / (\mathbf{1 \text{ month}})$



3. AOD synchronization: methodology

- When performing re-reconstruction, each Tier-1 produces the RECO and AOD format corresponding to the RAW data it has in its custody and exports it to all other Tier-1s.
- The full export to all sites should not take more than 2 weeks.
- For each exporting period of 2 weeks, the bandwidth from the exporting site is computed by:
$$(6 \cdot \text{fraction of AOD} / (2 \text{ week-period}))$$
- Proportional to full data sample size



4. T1s receiving MCPROD: methodology

- **All year around T2s to produce a number of MC events of 2030 Mio Events**
- **Half of the CPU are devoted for MC production, time for producing simulated event**
- **Volume aggregated at T1s according to child-parent affiliation**
- **Proportional to full data sample size and to the cpus size**



5. Simulating data analysis: methodology

- Each T1s has a share of RECO and AOD coming from 1st re-reconstruction
- Each T1 will receive requests for data skimming from a share of 1000 CMS Users according to its size
- The analysis pattern is assumed to be:
 - Half of the users will perform twice per month a “standard” analysis (i.e. importing 1% data)
 - 10% of the users will perform once per month a large dataset analysis (importing 10% data)

data set analysis being re-RECO/AOD (from any Tier 1) and full AOD (from parent site)

- Data rate evaluated as daily export set/1day
- Linear with the full event size and the number of active users (1000 in this calculation)



The results

- 1 slide per regional Tier1
- Pledged Cores are given in the old kSI2k units
- Remember: this are really raw values
- Links:
 - Solid line: sustained bandwidth (**data taking ONLY**)
 - Broken line: peak bandwidth (**may happen at any time: numbers shown is the total if it all happens at the same time**)
- For each Tier 1, the fraction of served users for analysis is a combination based on
 - Relative size child T2s for analyzing the share of 1srt AOD at considered Tier1
 - Relative size of T1 for analyzing the full AOD

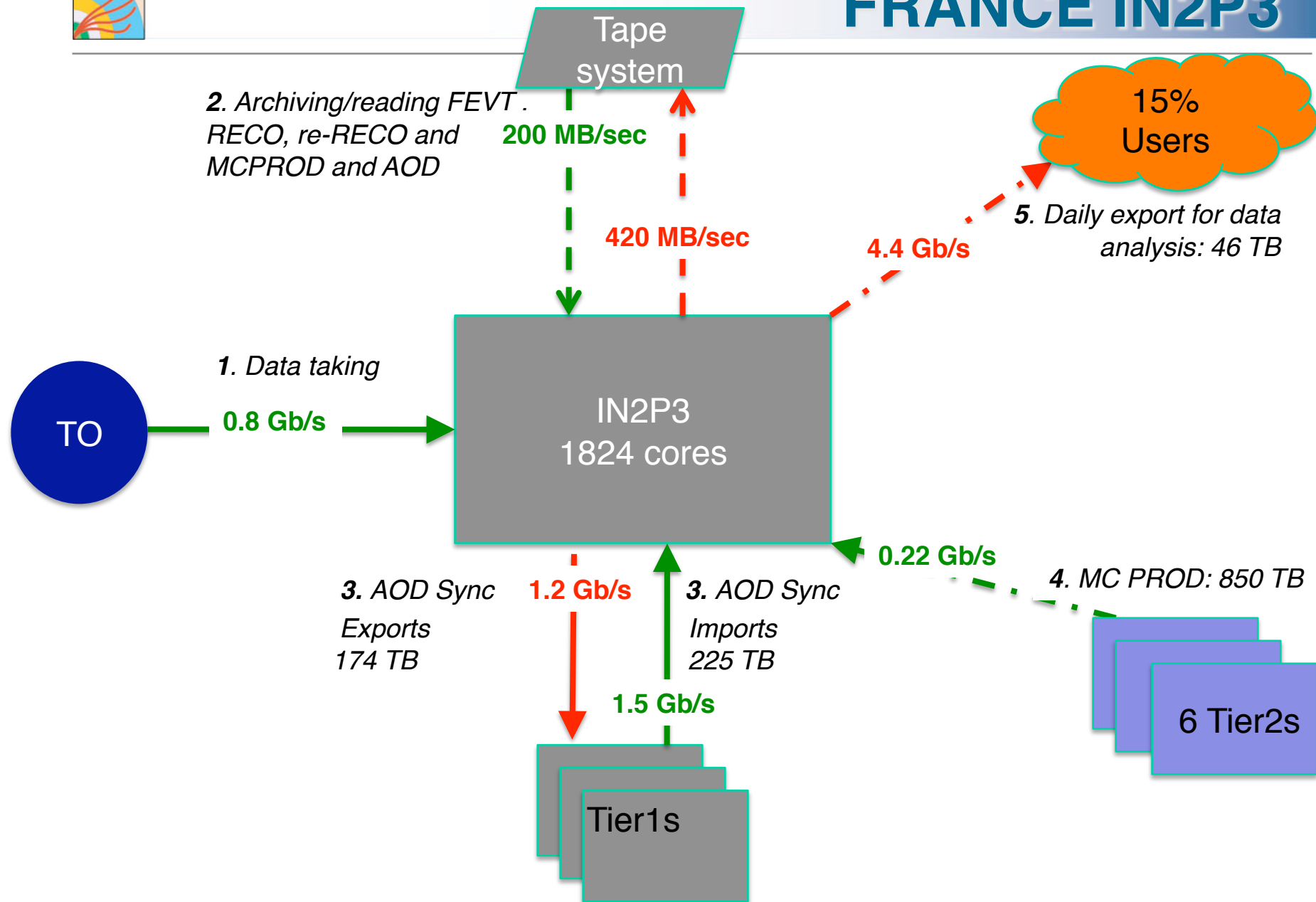


Next steps

- **We are discussing with each Tier 1 individually to see how far apart we are**
- **We need to refine the model**
- **These two steps need to be done before the CMS Compute Resource Board of June 18th**
 - **Final green light for publication of the 2009-2010 edition**
- **For the future: to be reviewed at the end of each data taking period**

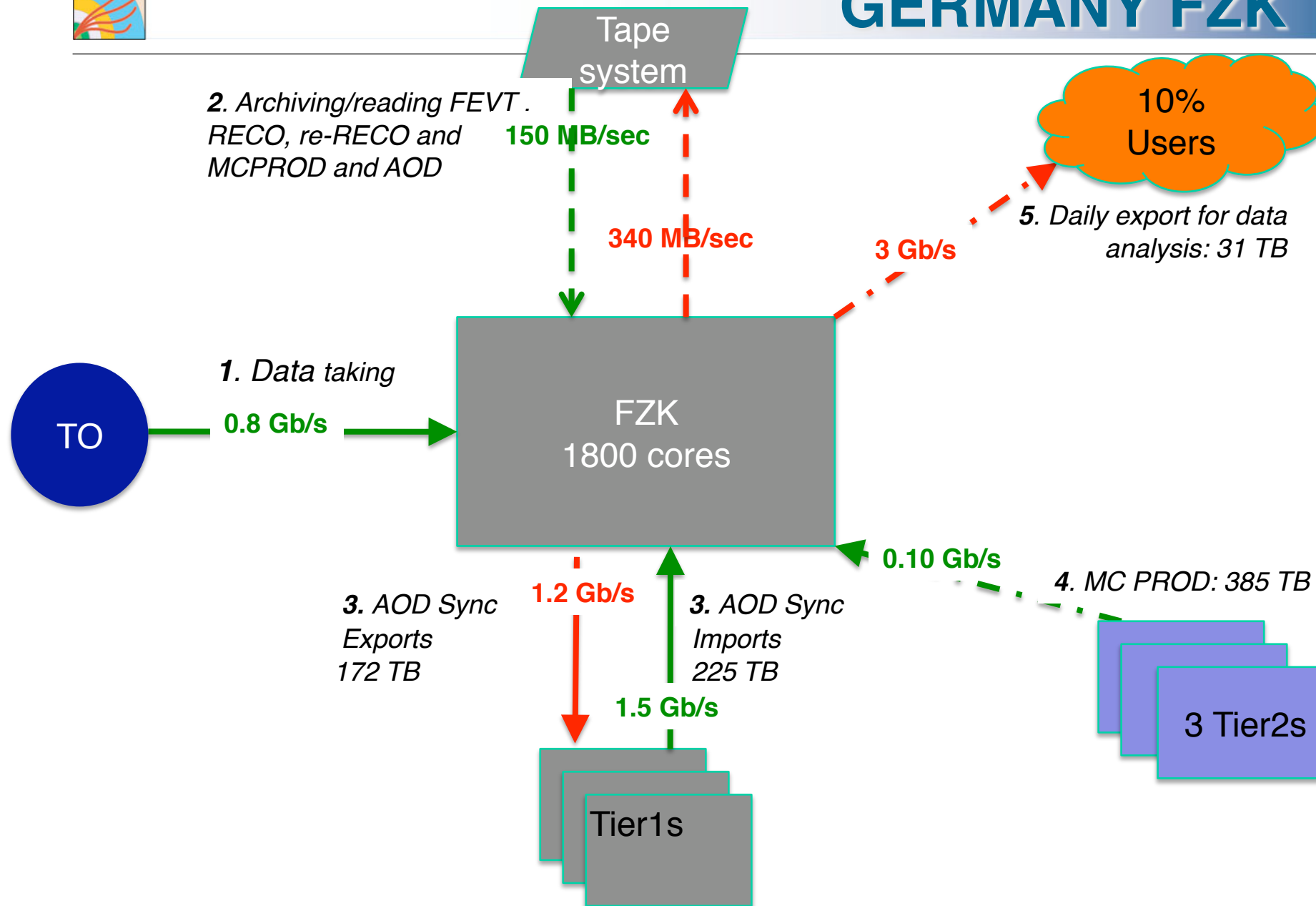


FRANCE IN2P3



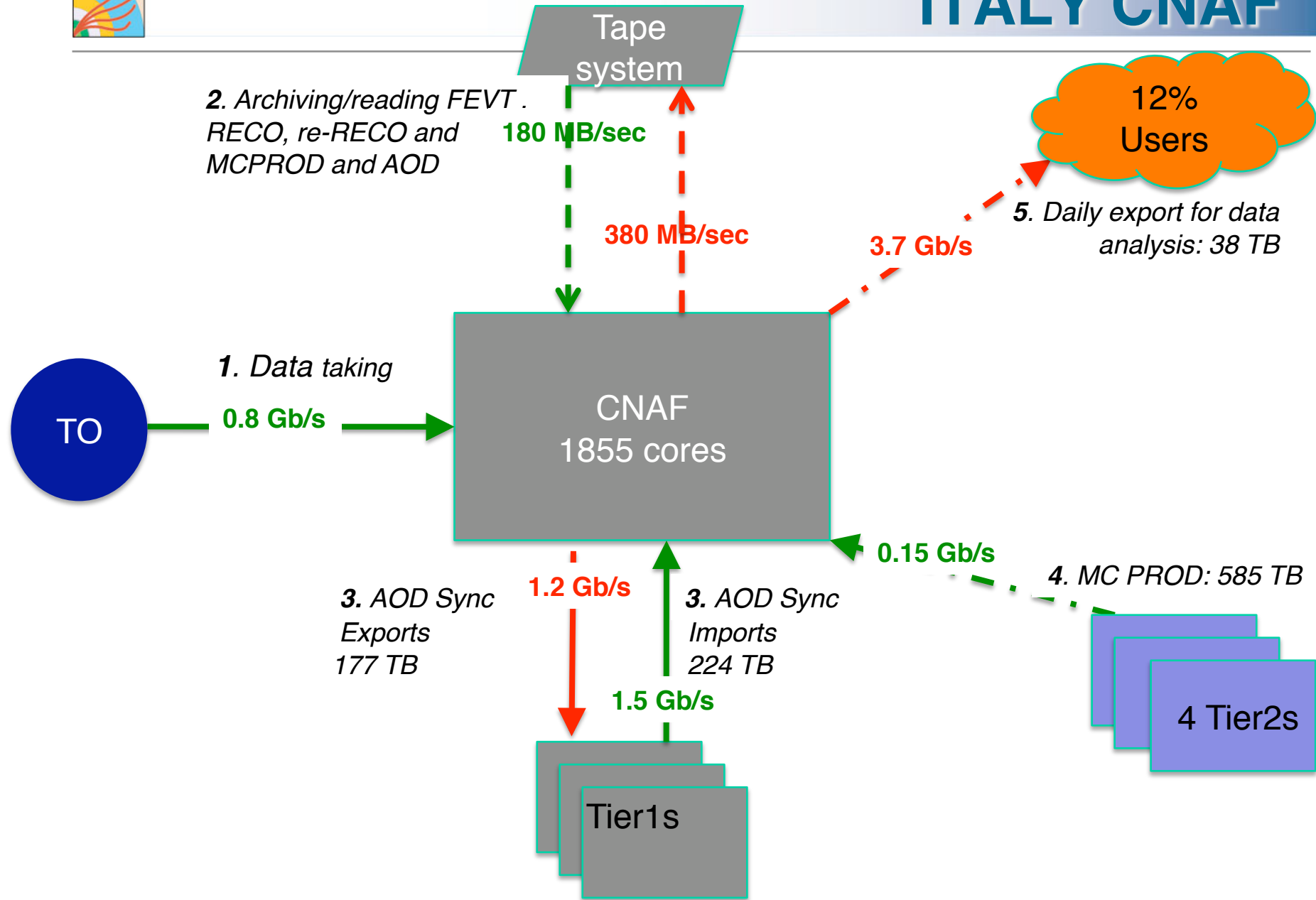


GERMANY FZK



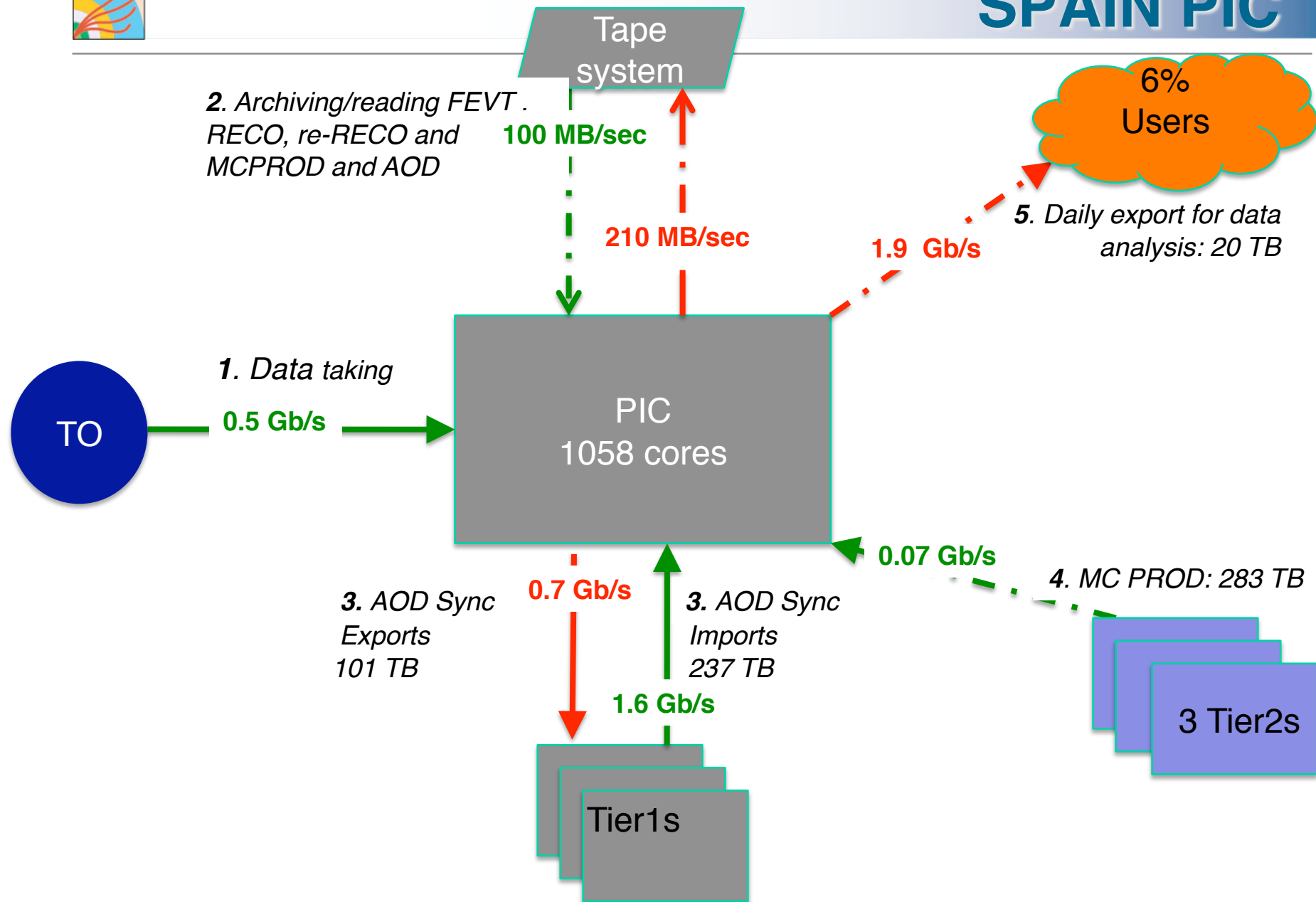


ITALY CNAF



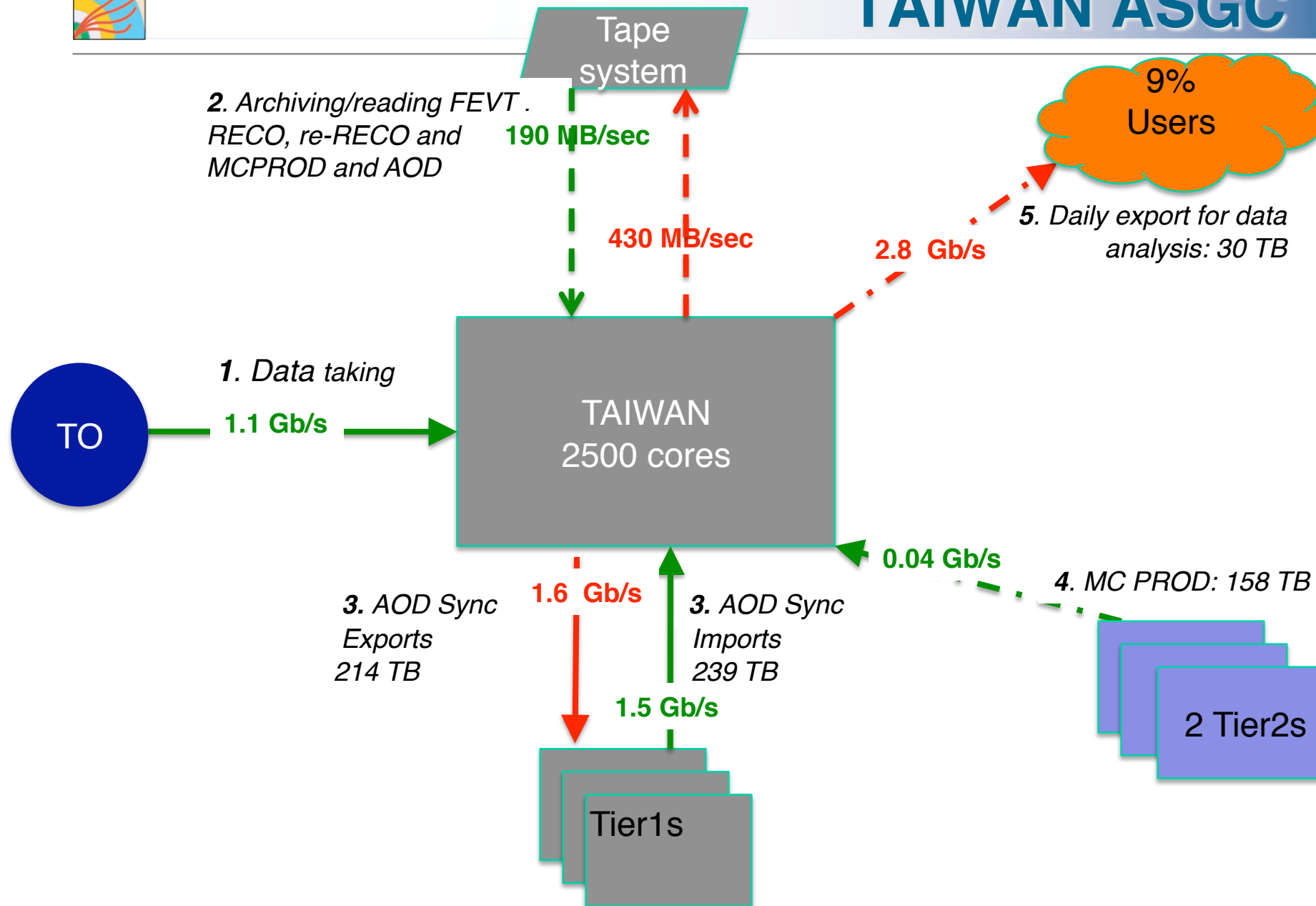


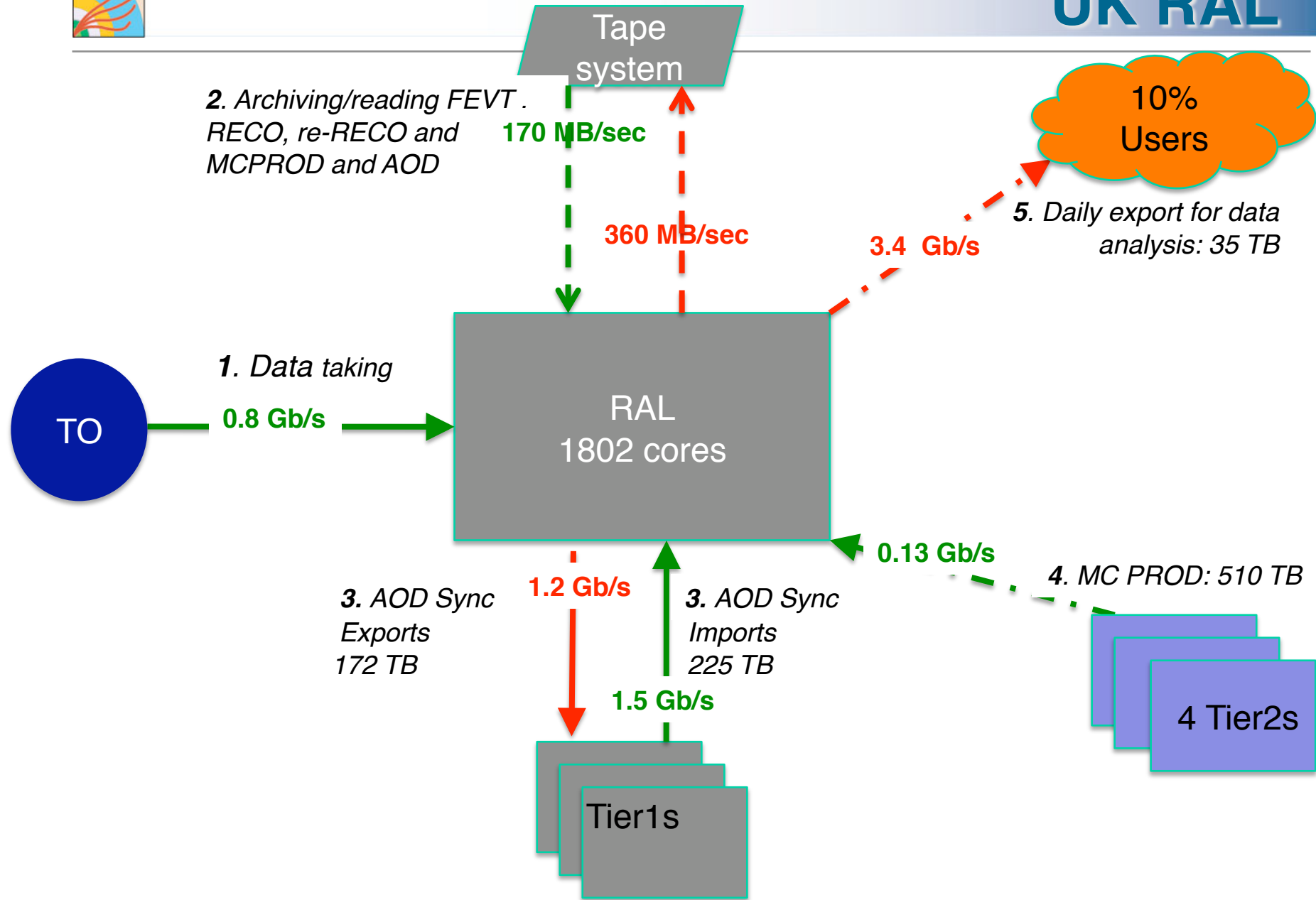
SPAIN PIC





TAIWAN ASGC







USA FNAL

