

Optimise resource usage and operation of lightweight Grid sites

S. Jézequel

ADC weekly – 29 Jan 2019

Introduction

- Overview : Recommendation (ICB-2018) to redirect funding from storage to CPUs for lightweight Grid site (2018 limit : 460 TB, 2019 (+15%): 520 TB)
 - Already implemented on voluntary basis on 5 sites : also had small amount of CPUs (<200 cores) (slide 3)
- Target : Propose a simple and reliable solution for bigger Grid sites (200-1000 cores)
 - Grid site : Host client commands to copy/read files on remote Grid SE
- Diskless site with direct access to associated SE : Current simplest solution compared to any change (xcache, ARC-CE)
- Diskless would require to avoid (BHAM expérience)
 - Network saturation
 - Too much load on remote SE
 - Keeping high IO job would require a tool to limit their amount (not existing yet)
 - Proposition : Run only low IO jobs
- Question : Will it reduce significantly the Grid CPU resources for large IO activity ?
 - Would not accept reduction by more than few % (total Grid cores : 300k)

Current diskless sites

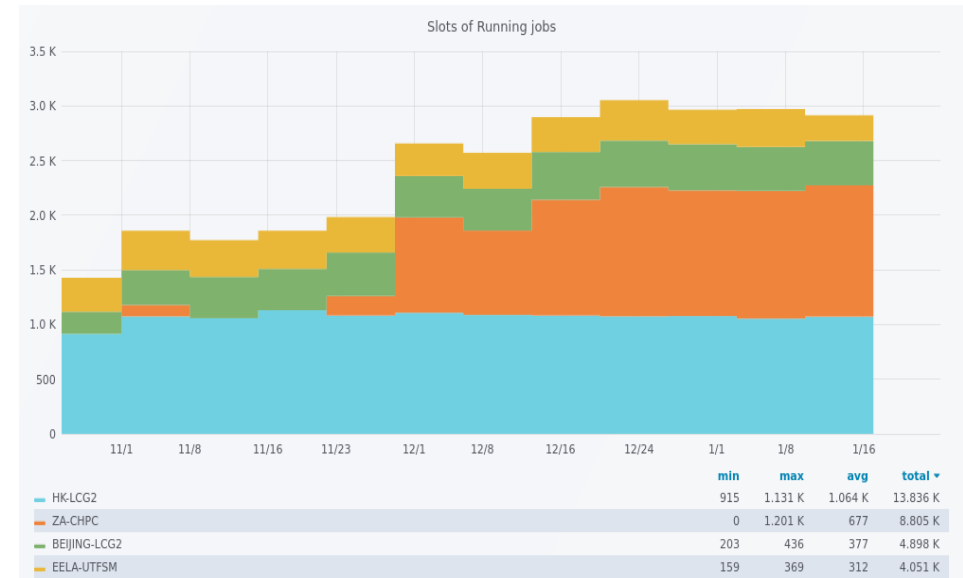


600 cores → 0.2 % of Grid capacity

- Running only low IO jobs → Sites were kept busy over last 3 months

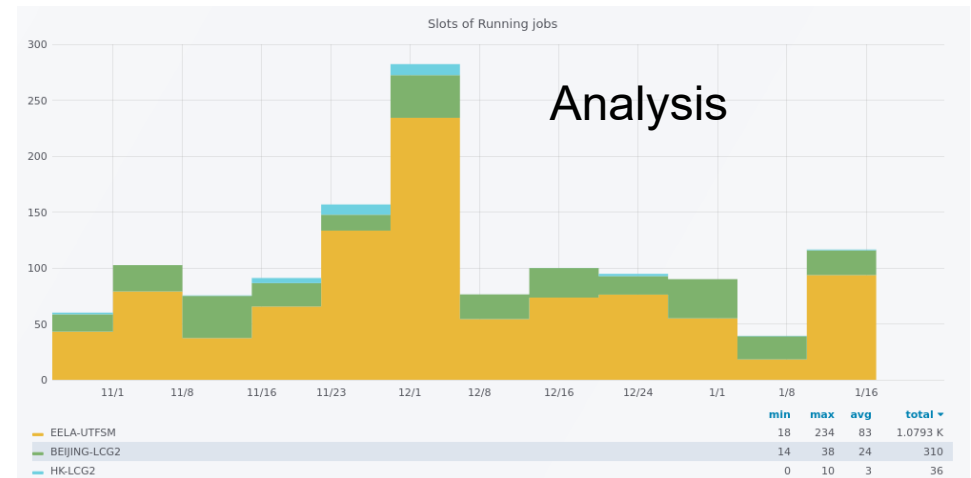
Isolated T2/T3 sites

Site	DATA DISK (TB)	SCRATCH DISK (TB)
HK-LCG2	420	22
ZA-CHPC	-	-
BEIJING-LCG2	310	60
EELA-UTFSM	360	30



3000 cores → 1 % of Grid capacity

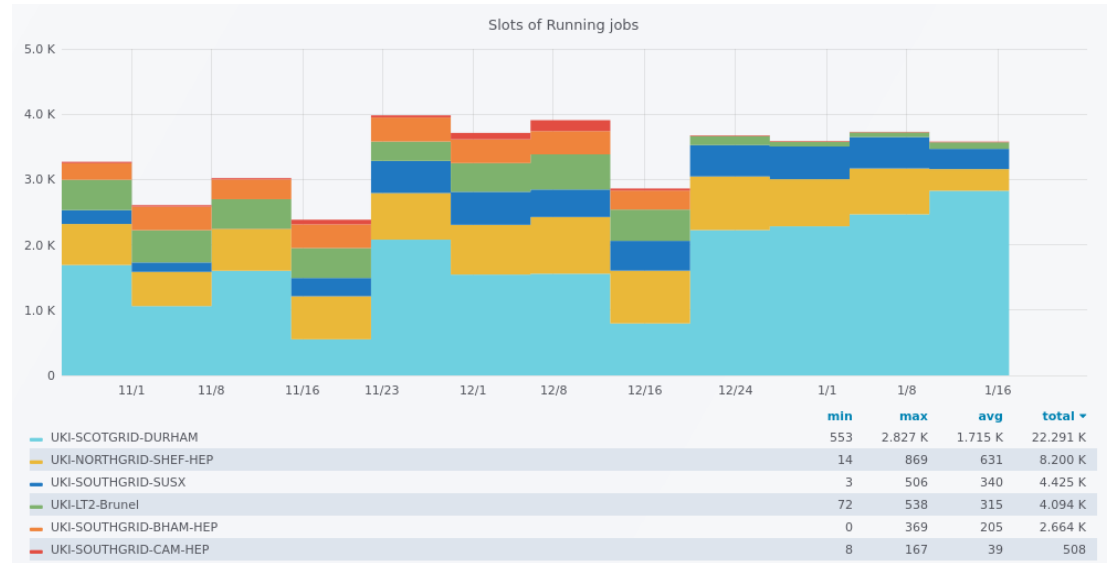
Storage capacity : ~1.2 PB



- ZA-CHPC (South Africa) (associated to INFN-T1) experience demonstrates no need for local DATADISK if low IO activity

Sites within T2 federation

Site	DATA DISK (TB)	SCRATCH DISK (TB)
DURHAM	192	-
SHEF	450	11
SUSX	10	-
Brunel	33	-
BHAM	-	-
CAM	190	8



Site	DATA DISK (TB)	SCRATCH DISK (TB)
WEIZMANN*	192	40
NCG-INGRID-PT*	177	11
PSNC	341	22
RO-02-NIPNE*	319	16



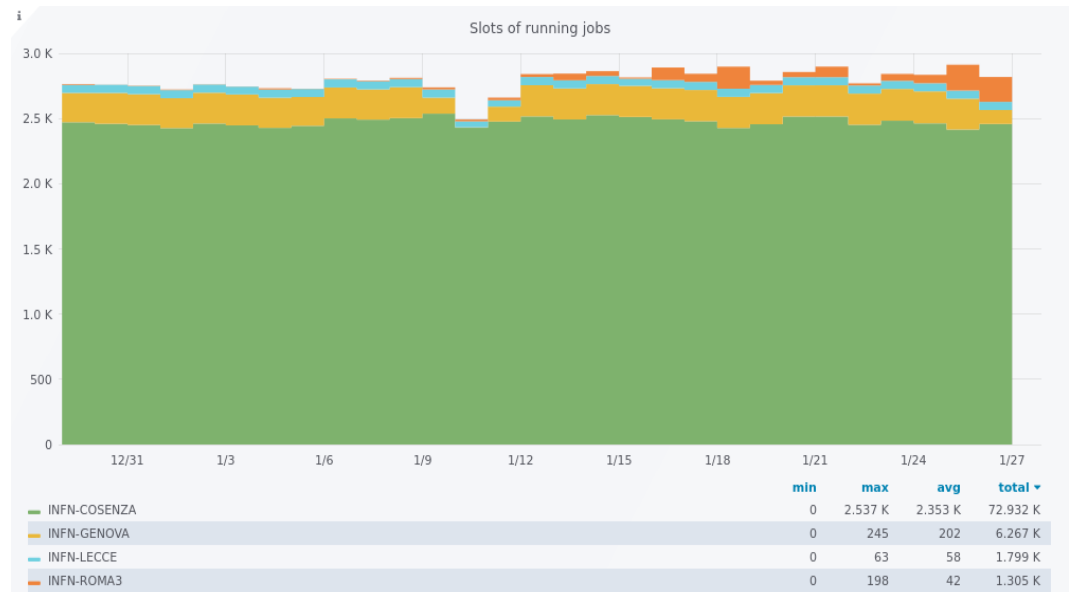
8000 cores → 2.5 % of Grid capacity

Storage capacity : ~2 PB

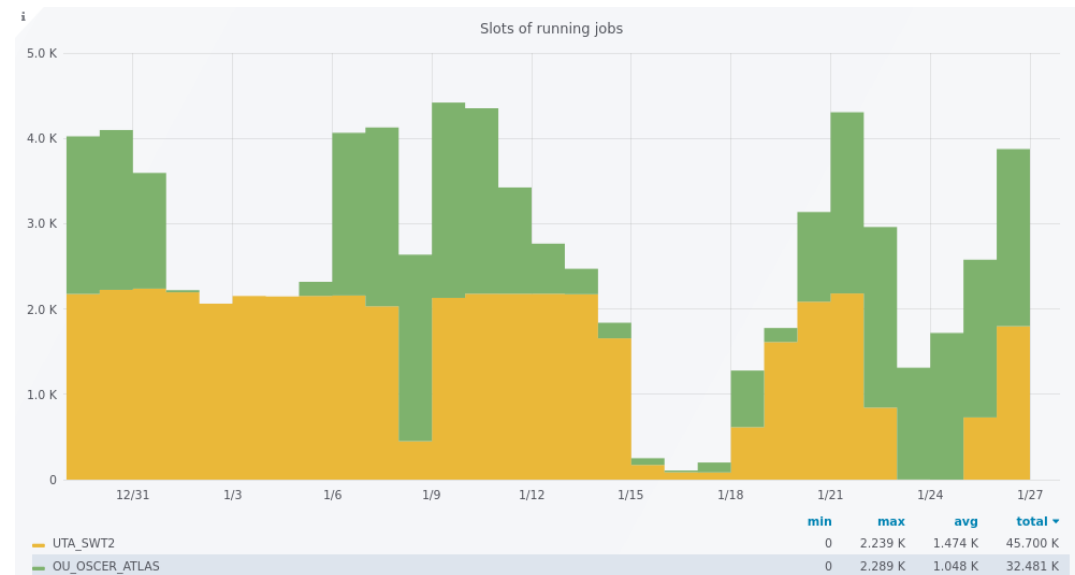
* : Site broken for analysis > 2 months

Sites within T2 federation (2)

Site	DATA DISK (TB)	SCRATCH DISK (TB)
COSENZA	330	55
GENOVA	1.1	4.4
LECCE	1.00	5
ROMA3	3.3	7.7
Bologna-T3	3	5



Site	DATA DISK (TB)	SCRATCH DISK (TB)
NERSC	192	-
UTA-SWT2	325	-
OU_OSCER_ATLAS	400	50

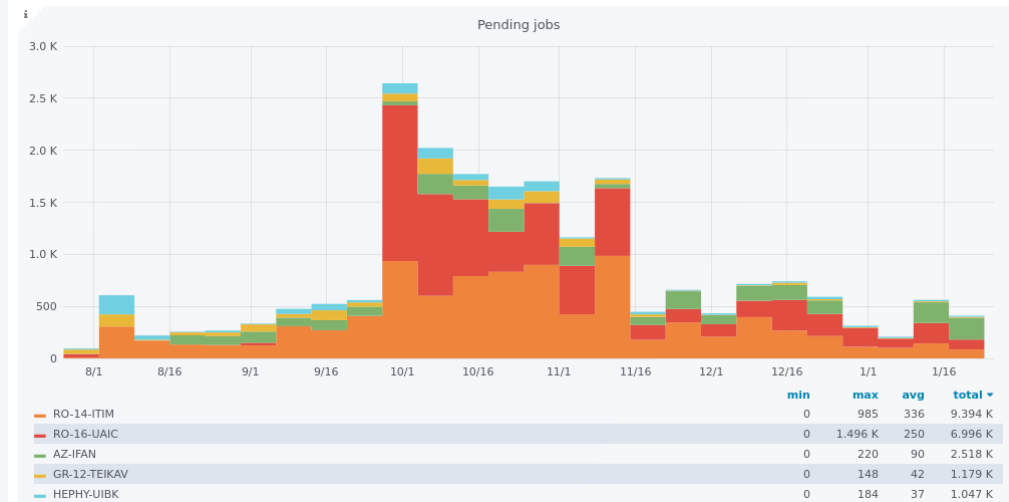
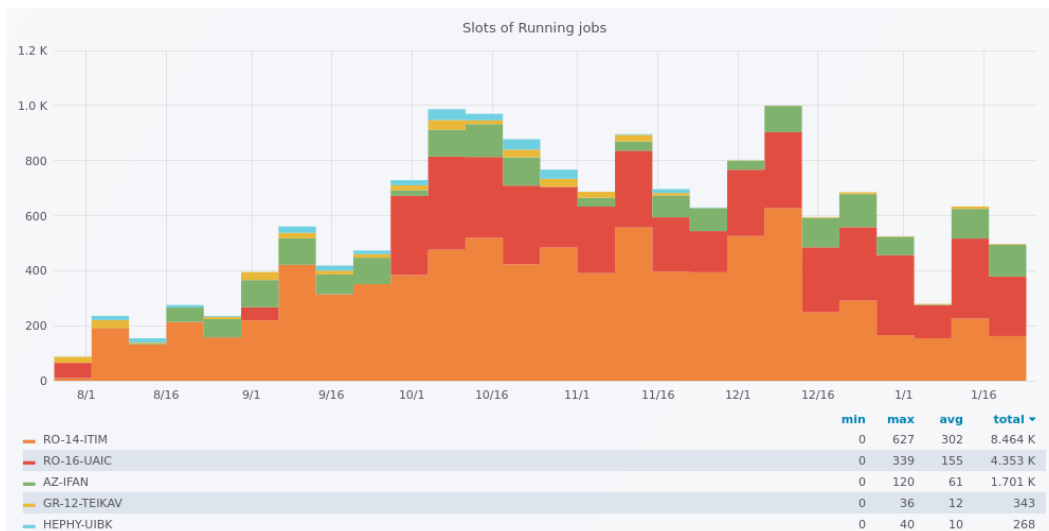


6000 cores → 2 % of Grid capacity

Storage capacity : ~1.4 PB

Enough activity with simul/evgen ?

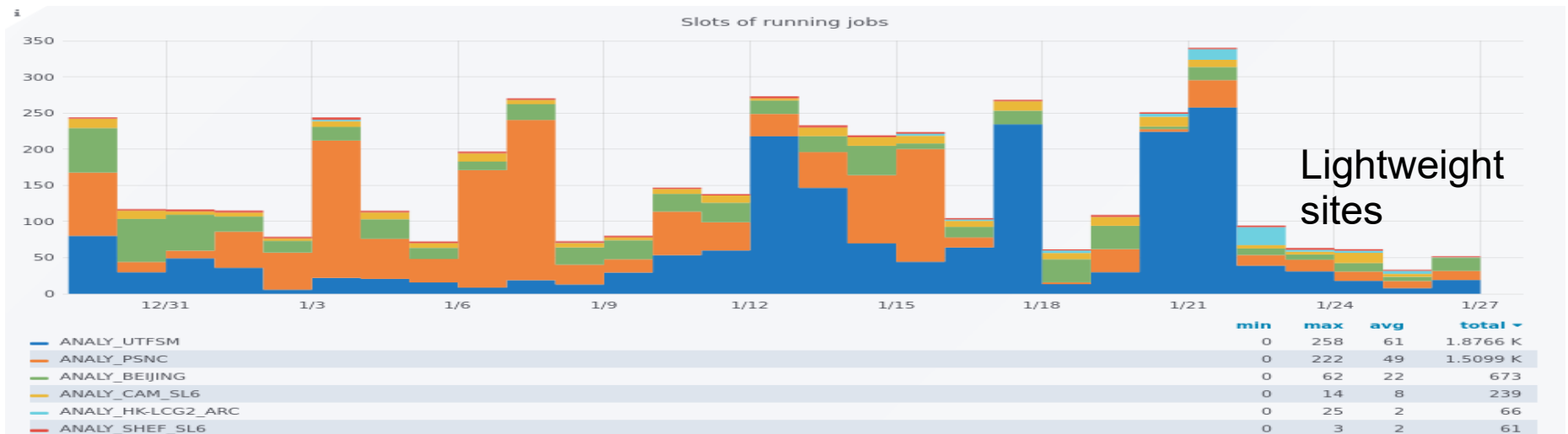
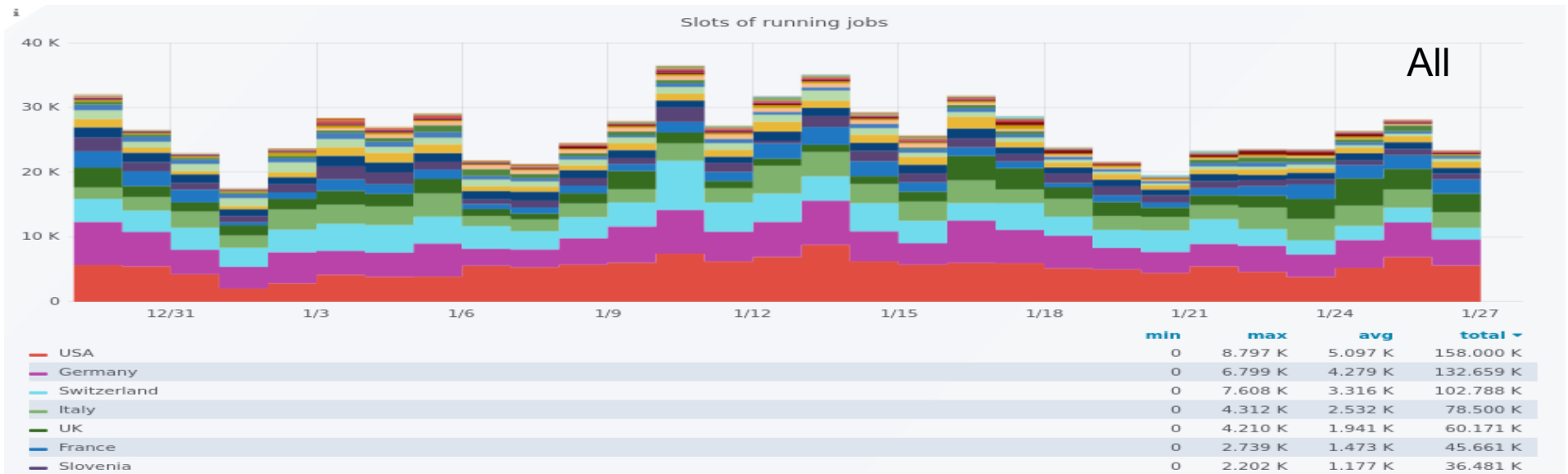
- Simul and evgen processed by
 - HPC
 - Cloud computing
 - Volunteer computing
- Reduction of CPU needs with increased usage fast sim ?
- Enough request to use Grid diskless sites in 2019 ? Beyond ?



→ Current diskless sites kept busy in last 6 months with simul only

Loose CPU capacity for analysis ?

No storage → No more analysis accessing local storage ?



→ Would loose O (1%) of analyse capacity

Recommendation for lightweight sites

- Minimise changes requested to the site
- Focus activity to low IO production jobs
 - Small Grid site represents ~ 5 % of Grid capacity
 - No damage to avoid to run high IO in these sites
 - Low IO → Almost not affected by network occupancy → Lower operational burden
 - No more analysis queue (currently~2%)→No interest to keep DATADISK
 - **Migrating to ATLAS@Home would be lowest support solution (still monitoring issue ?)**
- If site wants to keep storage for some time, reallocate all Grid storage to SCRATCHDISK and contribute to storage of second replicas of user outputs (lifetime of 2 weeks → storage decommissioning could be fast)
- Need reliable solution to avoid GRID storage for LOCALGROUPDISK (with remote site storage management)
- No more need to request local squid ?
- New switcher is mandatory (switcher3 to be validated and deployed soon by Jose)

Recommendation for lightweight sites (2)

- In 2019, transform lightweight sites to run only low I/O production jobs
 - Looses only ~15k slots for production
 - Only request to site admin : Close local analysis queue
- Beyond 2019 : since more sites are expected to become diskless, necessary to setup ATLAS site configuration to enable reasonable level of high I/O jobs according to network connectivity
 - Stay diskless : All informations are available (network connectivity per experiment, data flow rate per job type) but requires implementation in the workflow
 - Use cache if usefull for production (ARC-CE, xcache,...)
 - Could also be usefull for commercial cloud or HPC
- Up to the cloud to decide and organise the new shares
 - Report from each cloud/country on their decision by end february (before jamboree)