



ENABLING ATLAS BIG DATA PROCESSING ON PIZ DAINT AT CSCS

Gianfranco Sciacca

AEC - Laboratory for High Energy Physics, University of Bern, Switzerland

CHEP Conference 2019 - 4-8 October 2019, Adelaide, Australia

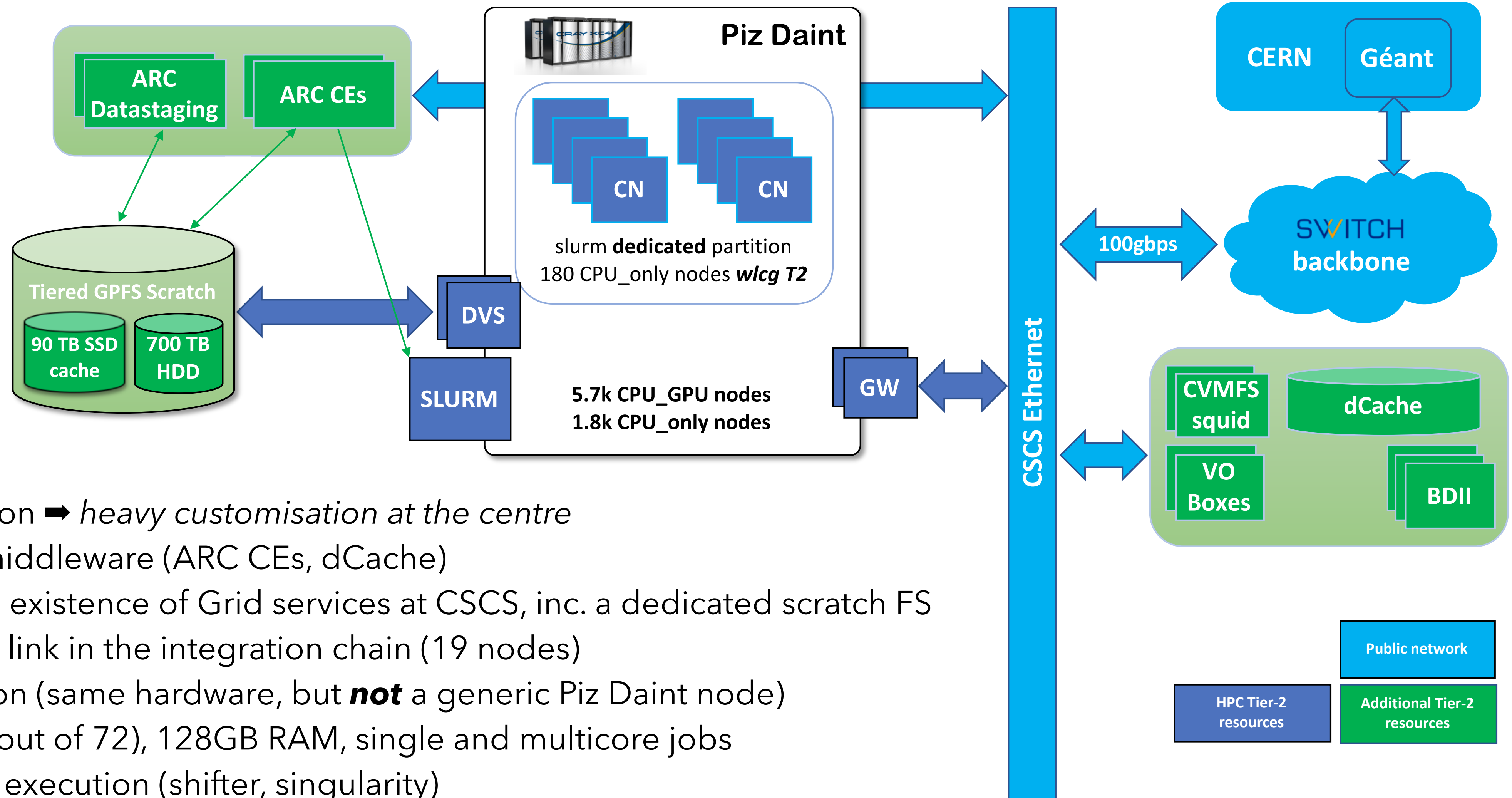
u^b

- ▶ **HPC integration with the LHC experiment frameworks has been a hot topic for several years**
- ▶ **Hoping for:**
 - ▶ more computing for the same price
 - ▶ access to high end technologies
 - ▶ reducing operational costs
- ▶ **Is it really so?**

The Swiss HEP computing community and CSCS started working on the HPC integration with the LHC experiment Tier-2 facilities in 2014

- ▶ **CSCS has hosted a WLCG Tier-2 site for ATLAS, CMS, LHCb for ~10 years** - Dedicated commodity cluster (*Phoenix*)
- ▶ **2014/15 - Cray XK7 integration, pioneered by ATLAS**
 - ATLAS Geant4 in production for 6 months
 - Remote submission to CSCS with a modified ARC CE
- ▶ **2015 - early tests and PoC on a Development Cray XC 40**
- ▶ **2016 - last procurement of dedicated Tier-2 hardware for Phoenix**
- ▶ **2017 - started Tier-2 production on Piz Daint (1.6k cores) alongside Phoenix**
 - Tackled most of the integration issues
 - Integrated **all** experiment workflows
- ▶ **2017 - ATLAS scale up test: 27k cores with Geant4**
- ▶ **2017 - decided to phase out dedicated resources → expect +27% CPU year-on-year**
- ▶ **2018 - ATLAS (followed by CMS) evaluated Tier-0 reconstruction "on demand"**
 - Hardening of the HPC side for full-scale operations
 - 12k cores, +1 PB storage
- ▶ **2019 Tier-2 facilities fully migrated to Piz Daint (~12k cores in production), Phoenix phased out**

TIER-2 ON PIZ DAINT, FULLY MIGRATED BY 1/3/2019



Highlights:

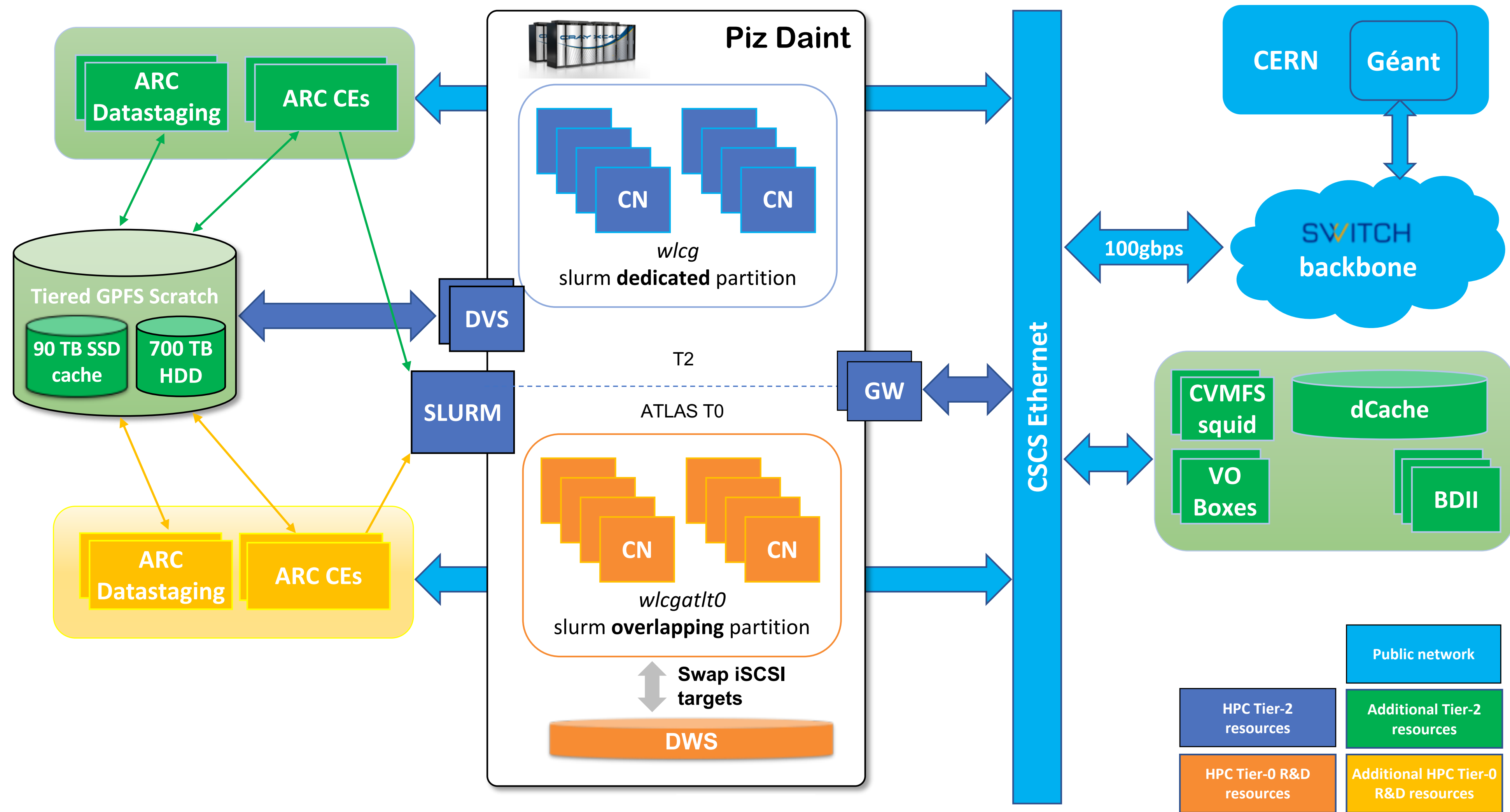
- Full Tier2 integration → heavy customisation at the centre
- Standard WLCG middleware (ARC CEs, dCache)
- Eased by previous existence of Grid services at CSCS, inc. a dedicated scratch FS
- Cray DVS a crucial link in the integration chain (19 nodes)
- Node customisation (same hardware, but **not** a generic Piz Daint node)
- 64 threads/node (out of 72), 128GB RAM, single and multicore jobs
- Containerised job execution (shifter, singularity)

u^b

- ▶ **Getting this far took 2 years (+1.5y for preliminary projects)**
 - ▶ several policy relaxed from the HPC side
- ▶ **Solutions coupled to the specific machine environment**
 - ▶ and to the machine lifetime
 - ▶ all/most/some of this might change with the next generation HPC

- ▶ **ATLAS RAW data reconstruction is performed at CERN on dedicated “hardened” resources**
 - ▶ Memory hungry and I/O intensive, exceeding the Piz Daint specs
 - ▶ ATLAS used spill-over to Grid sites in late 2018 to process some PbPb heavy ion runs
- ▶ **Integration very laborious (April-July 2018)**
 - ▶ Based on the Tier-2 architecture, but nodes no longer reserved (re-configured on the fly)
 - ▶ Much extra tuning needed, to address memory shortage and I/O load
 - ▶ E.g. in slurm, DVS, DataWarp, etc.
- ▶ **All in all promising results**
 - ▶ ATLAS *Physics_BphysLS* stream reconstruction validated (~10% of full run)
 - ▶ **But:** for a full run reconstruction: **50% of each node CPUs left idling**
 - ▶ Hampered by bugs uncovered in the CRAY *DVS* and *Data Warp* layers
 - ▶ Bugfixes for CRAY *DVS* and *Data Warp* should have recovered that

PIZ DAINT WITH TIER-0 R&D EXTENSION



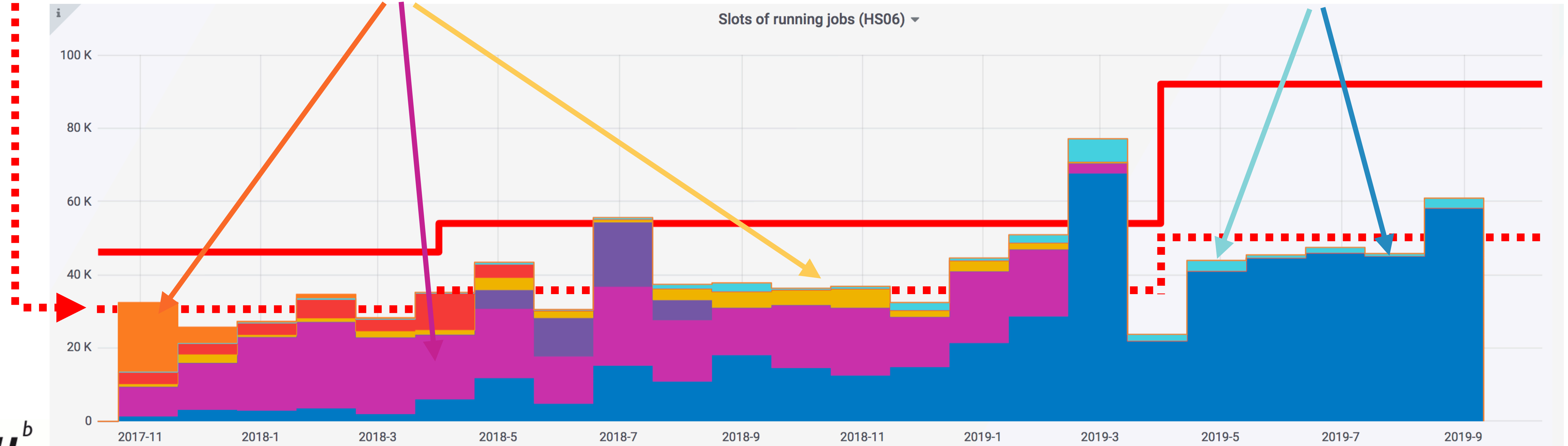
TIER-2 PERFORMANCE – HS06 DELIVERY

ATLAS 2-YEAR VIEW

Pledge

Phoenix

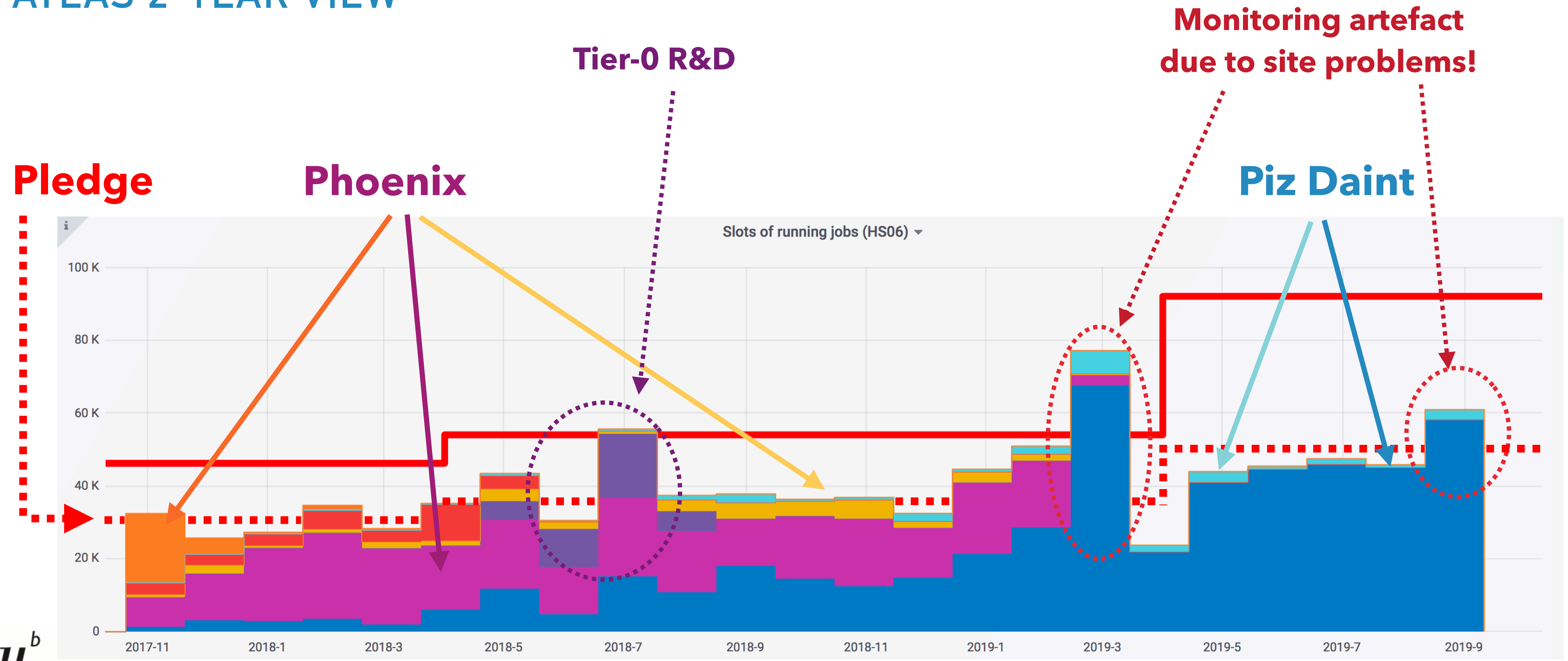
Piz Daint



u^b

TIER-2 PERFORMANCE – HS06 DELIVERY

ATLAS 2-YEAR VIEW



u^b

TIER-2 PERFORMANCE – EFFICIENCIES (ATLAS)

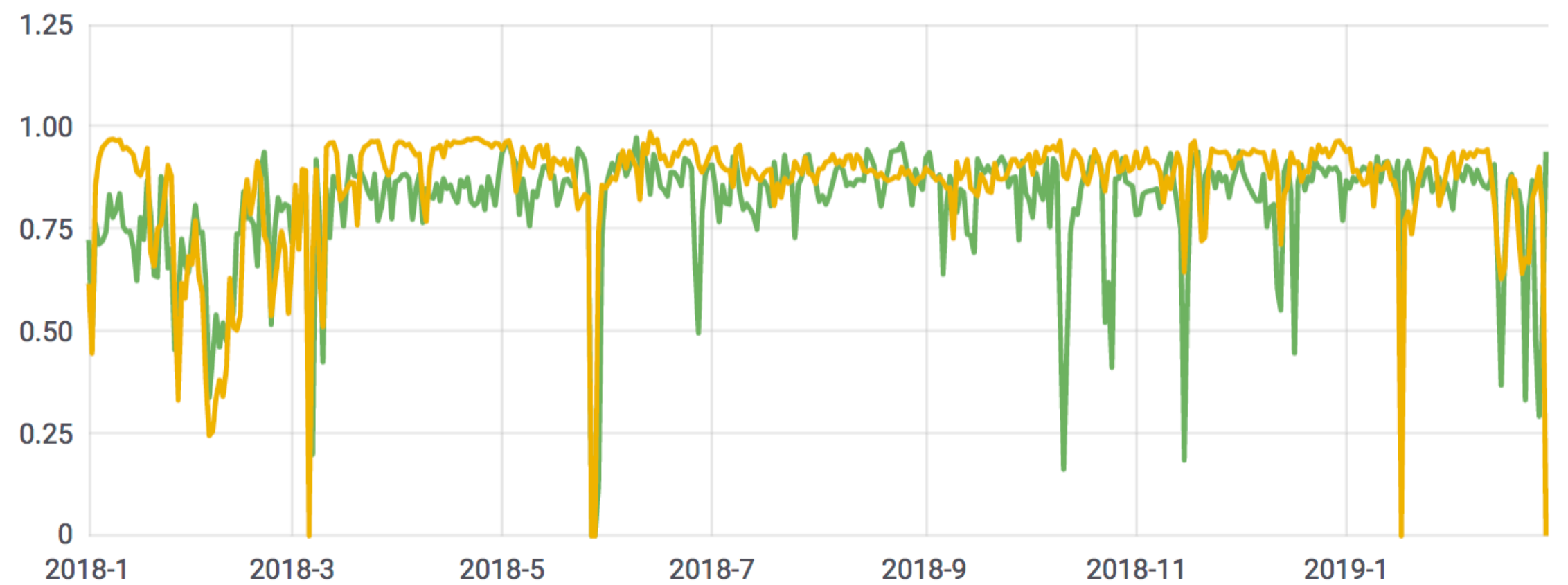
WC of successful vs failed jobs



CPU / WC efficiency of good jobs



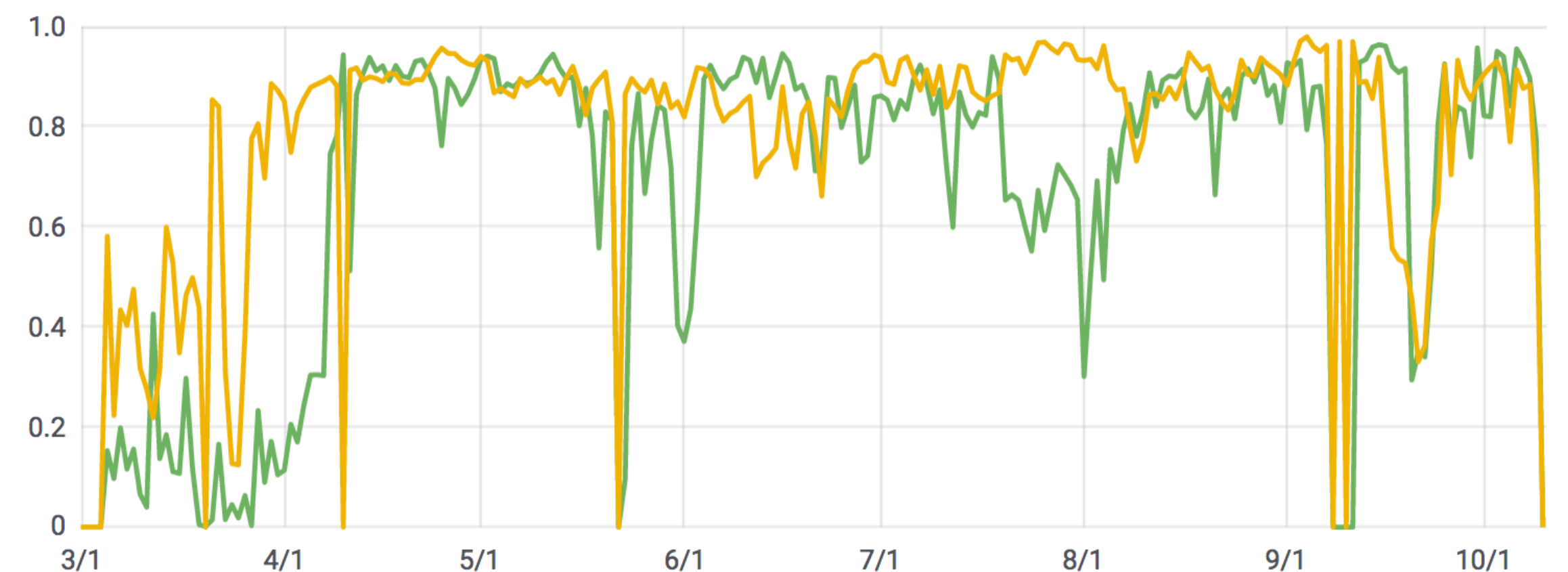
CPU Efficiency: Good jobs ▾



| | min | max | avg | current |
|-----------------|-----|-------|-------|---------|
| ANALY_CSCS | 0 | 0.973 | 0.811 | 0.939 |
| CSCS-LCG2_MCORE | 0 | 0.986 | 0.857 | 0 |

Phoenix

CPU Efficiency: Good jobs ▾



| | min | max | avg | current |
|---------------------|-----|-------|-------|---------|
| CSCS-LCG2-HPC_MCORE | 0 | 0.979 | 0.784 | 0 |
| ANALY_CSCS-HPC | 0 | 0.963 | 0.675 | 0 |

Piz Daint

▶ What is “High Performance” ?

▶ Integration costs very *high* (increased complexity vs a traditional cluster)

- ▶ Solution (likely to be) limited to the lifetime of the specific machine
- ▶ A Grid site running on nodes of a HPC (4-year old Xeon)

▶ Operational costs not decreased, rather *increased*

- ▶ Tier-2 workloads often exceed the hardware specs (e.g. RAM, I/O)
- ▶ Struggling to keep the allocated nodes full (< 70% usage on average), and unsolved fair-share challenges

▶ HS06 delivery poor and efficiency *comparable* to a traditional cluster

▶ Cost of allocation “artificial” (set by provider)

- ▶ Initially favourable (+27% CPU/year)
- ▶ Predicted to be 15% for the next HPC

▶ Loss of flexibility

- ▶ Committed to a machine for its lifetime
- ▶ Harder to adapt to computing model changes

- ▶ **Swiss HEP and ATLAS running a full WLCG Tier-2 on the HPC Piz Daint @ CSCS**
- ▶ **No gaugeable benefits, rather penalties. Costly**
 - ▶ Unless additional opportunistic usage, but not our case
- ▶ **The key is in the funding model**
 - ▶ Dedicated paid-for nodes cannot be filled (our case)
 - ▶ Assured node-hour allocation per funding cycle would make a difference (hard to negotiate)
 - ▶ HPC paid "by the node" will not help in the HL-LHC challenge
- ▶ **Accelerators are common on many clusters that are not commonly hyped as "HPC" or in TOP500**
 - ▶ Those are the ones who also agree on *managed* opportunistic usage ⇒ *_more_ computing for the same price*
- ▶ **Swiss ATLAS is re-considering running a Tier-2 site on a costly HPC**
 - ▶ Politically hard at the current hype level

THANK YOU FOR YOUR ATTENTION!

BACKUP

- ▶ **172 dedicated Cray Nodes, 64 cores/node (out of 72), 128GB RAM, single and multicore jobs**
 - ▶ OS/Software on the nodes is containerised (shifter, singularity)
- ▶ **CVMFS tiered cache (no local disks on the nodes)**
 - upper layer 6GB in RAM
 - lower layer preloaded on GPFS (scratch)
- ▶ **GPFS scratch filesystem (increased metadata performance over Lustre/CRAY Sonexion) with SSD cache layer**
- ▶ **Cray's DVS nodes needed for mounting ARC session, cache and CVMFS shared cache**
- ▶ **SLURM dedicated partition, submission from ARC CE**
 - Cores are consumable resource, memory is not
 - This allows jobs to use more than the base 2GB RAM/core (no swap)
 - Overall memory consumption is surprisingly balanced
 - Specific QoS and PriorityTier to make the backfill scheduler pass 'fast' over LHC jobs
- ▶ **ARC SLURM batch generator modified to create batch jobs that:**
 - Have the proper flags (QoS, partition, nice, etc.)
 - Run the payload within a shifter container or singularity