

- ▶ Motivation and challenges
- ▶ Architecture
- ▶ Production performance
- ▶ Assessment and outlook



# PRODUCTION EXPERIENCE AND PERFORMANCE FOR ATLAS DATA PROCESSING ON A CRAY XC-50 AT CSCS

SWISS NATIONAL SUPERCOMPUTING CENTRE



**Gianfranco Sciacca**

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

## WLCG computing on HPC systems

- ▶ **The HENP computing community will face several challenges with respect to the computing requirements for the next decade and beyond**
- ▶ **For the LHC community, requirements for the High-Luminosity LHC runs (2025-2034) are expected to be a factor of ~50 higher than today**
- ▶ Available **budgets** are expected to be **flat at best**
- ▶ Novel computing models making a more dynamic use of heterogeneous resources (supercomputers and clouds) need to be evaluated in order to address such challenges
- ▶ HPC machines are increasingly powerful, could play a crucial role in delivering **more computing for the same price**
- ▶ **One of the biggest challenges is the transparent integration with the complex experiment data processing frameworks**

WLCG computing on HPC systems

## ▶ HPC is awesome

### ▶ Piz Daint Cray XC50 / XC40 @ CSCS

- 'cpu/gpu hybrid' (5320 nodes) and 'multicore' (1431 nodes), 361,760 cores
- NVIDIA Tesla P100, Xeon E5-2690v3 2.6 GHz, 521 TB of RAM, 25k Tflops
- Cray Aries high-speed "dragonfly" topology interconnect
- Lustre Sonnexion 3000, 6.2 PB, Sonnexion 1600, 2.5 PB
- GPFS, 700 TB + 90 TB SSD transparent cache
- DVS, Burst Buffer, Data Warp (POSIX filesystem on demand on SSDs)

WLCG computing on HPC systems

## ▶ HPC is Awkward

### ▶ No local disk

- Breaks a lot of standard HEP Linux workflows

### ▶ Minimal OS

- Designed to accelerate parallel software
- Many expected Linux tools are missing
- Runs a stripped down version of SUSE, and doesn't upgrade often

### ▶ Limited RAM

- Most of the CPU only nodes with 1 GB / core, some with 2 GB / core
- No swap

### ▶ Network connectivity not guaranteed

- Must be negotiated
- Needs gateways
- Interfacing external services is not straightforward (e.g. mount a directory)

## WLCG computing on HPC systems

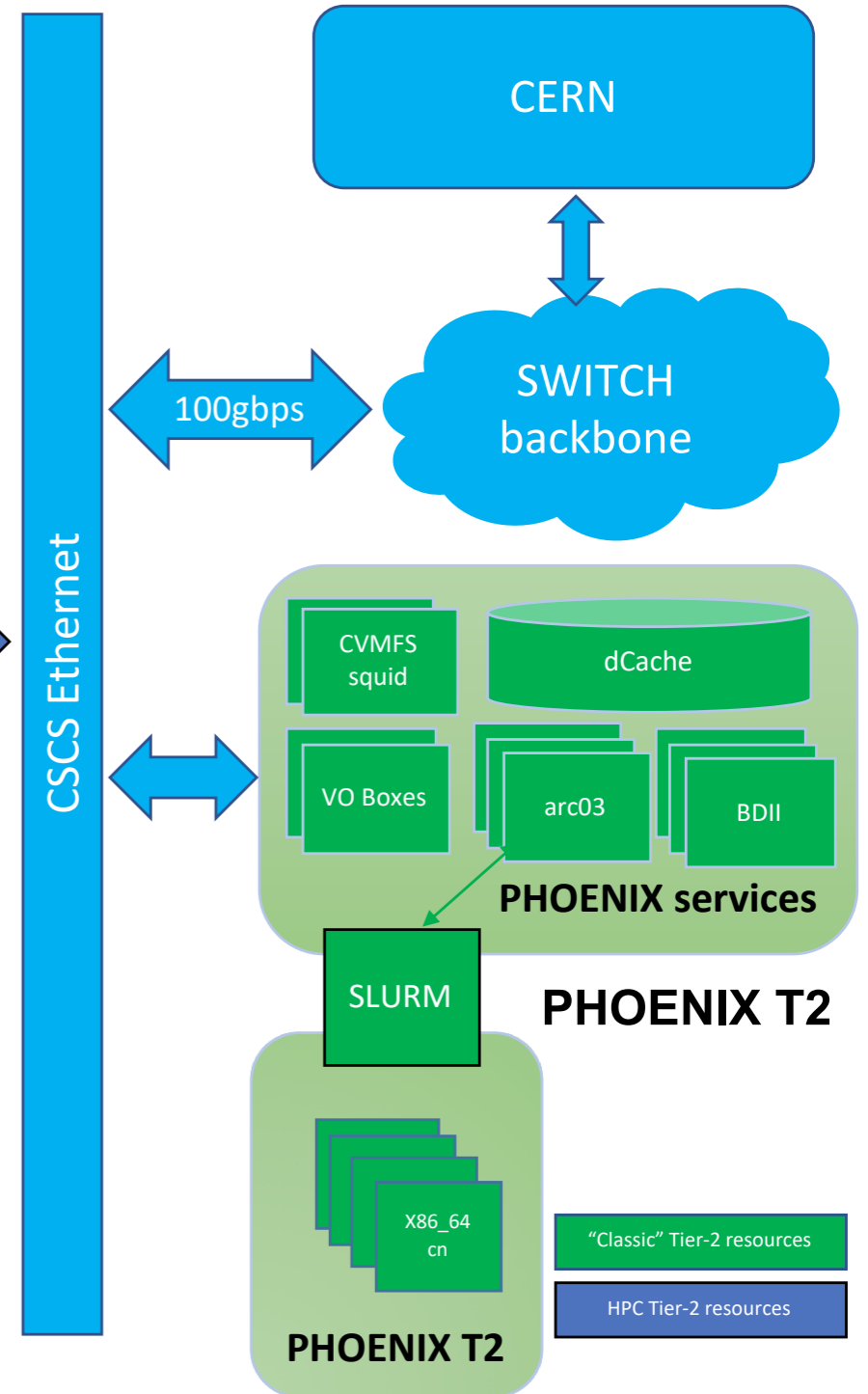
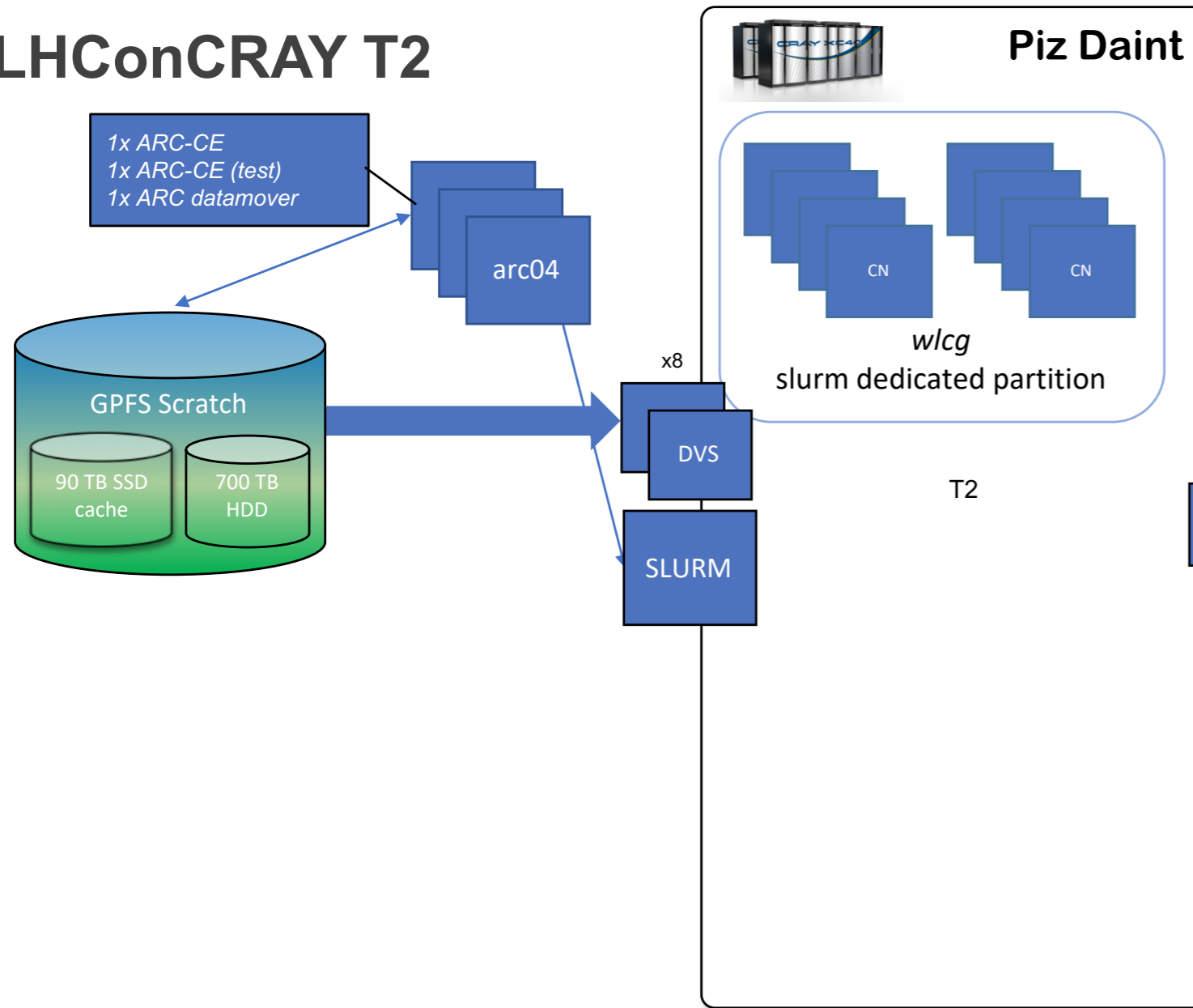
- ▶ **The Swiss HEP computing community and CSCS have started working on the HPC integration with the LHC experiment Tier-2 facilities in 2014**
- ▶ **ATLAS Geant4 simulation**
  - Ran in production for 6 months on a Cray XK7
  - Integrated by means of a modified ARC CE, submitting remotely to CSCS
- ▶ **LHConCray project (ATLAS, CMS, LHCb)**
  - Ran for about 2 years in 2016-17
  - Aimed at integrating Piz Daint with the LHC experiment frameworks
  - Targeted all experiment workflows (including user analysis)
  - Went in production with 1.6k cores in 2017
- ▶ **WLCG Tier-2 facilities migrated to Piz Daint**
  - Decision taken at the end of 2017
  - ~4k cores by April 2018, >10k by April 2019

Powered by



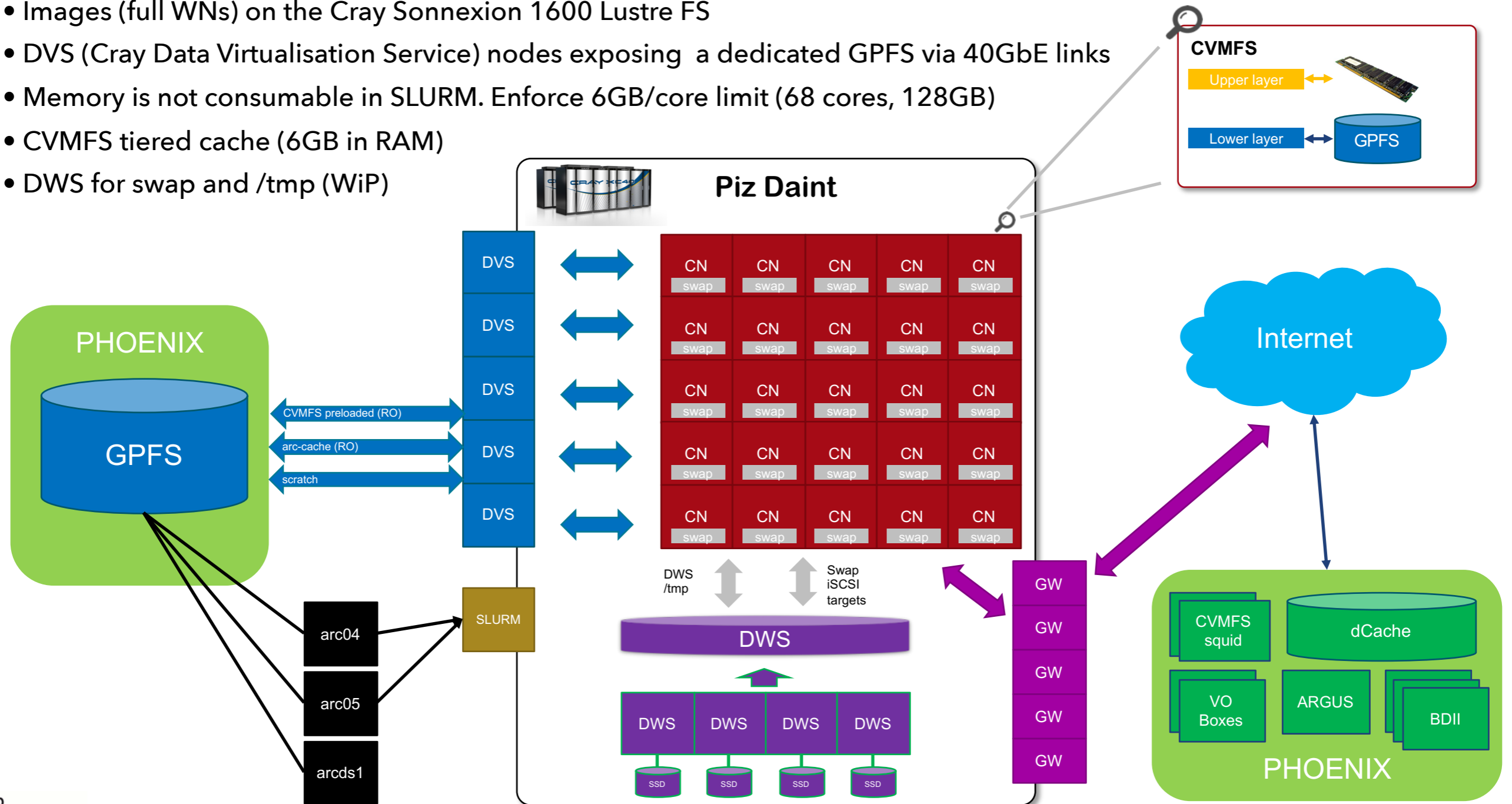
## Classic x86\_64 and HPC systems

### LHCOnCRAY T2



## HPC system highlights

- Jobs run in Docker containers using Shifter
- Images (full WNs) on the Cray Sonnexion 1600 Lustre FS
- DVS (Cray Data Virtualisation Service) nodes exposing a dedicated GPFS via 40GbE links
- Memory is not consumable in SLURM. Enforce 6GB/core limit (68 cores, 128GB)
- CVMFS tiered cache (6GB in RAM)
- DWS for swap and /tmp (WiP)



LHConCRAY - Acceptance Tests 2017 - Configuration evolution 1

## WLCG computing on HPC systems



### Slots of Running Jobs

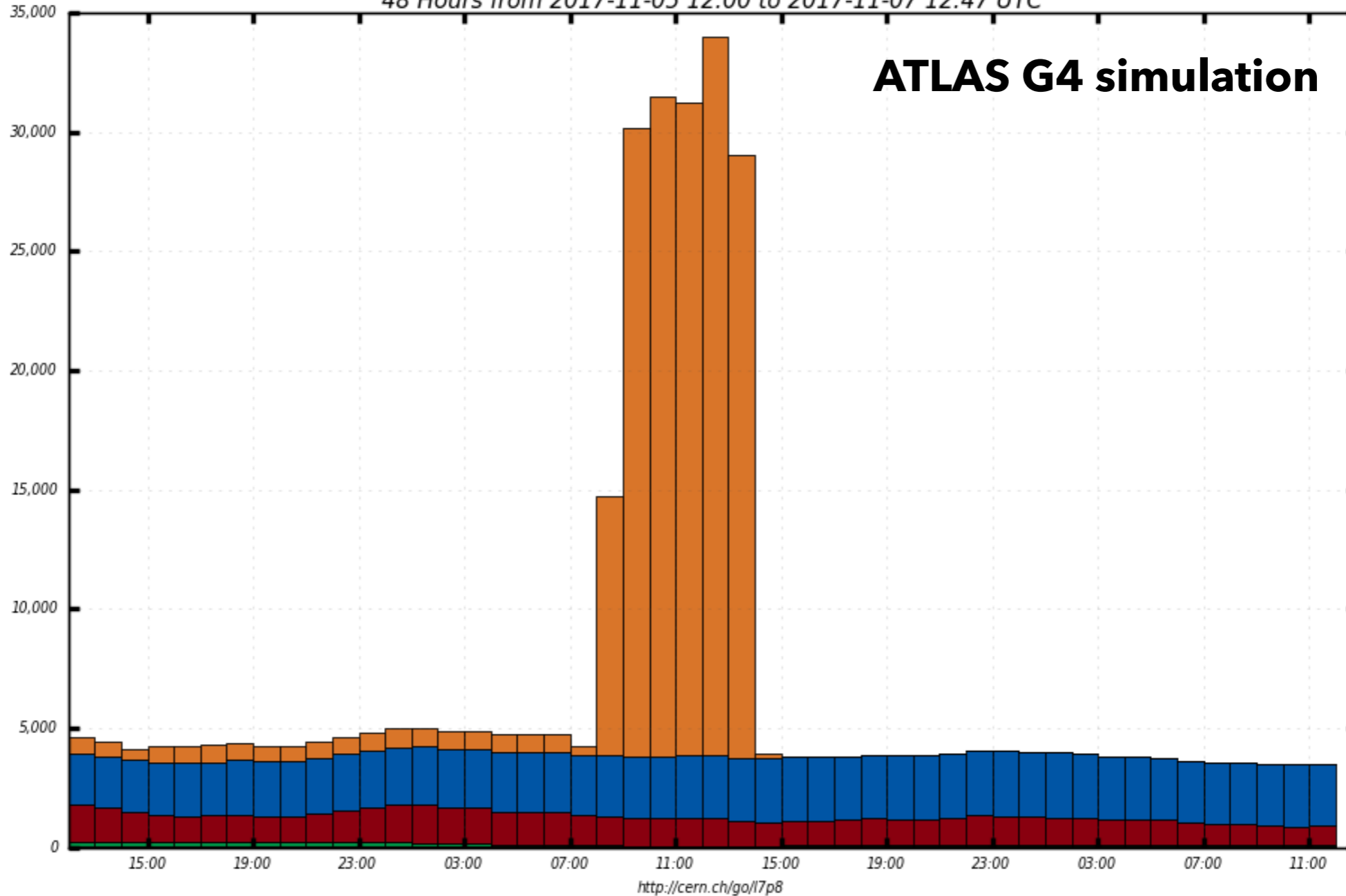
48 Hours from 2017-11-05 12:00 to 2017-11-07 12:47 UTC

## ATLAS G4 simulation

30.000



T2 slots



CSCS-LCG2-HPC\_MCORE\_TEST  
ANALY\_CSCS

CSCS-LCG2  
CSCS-LCG2-HPC\_MCORE

CSCS-LCG2\_MCORE  
ANALY\_CSCS-HPC

CSCS-LCG2-HPC

Maximum: 34,012 , Minimum: 3,466 , Average: 7,186 , Current: 3,478

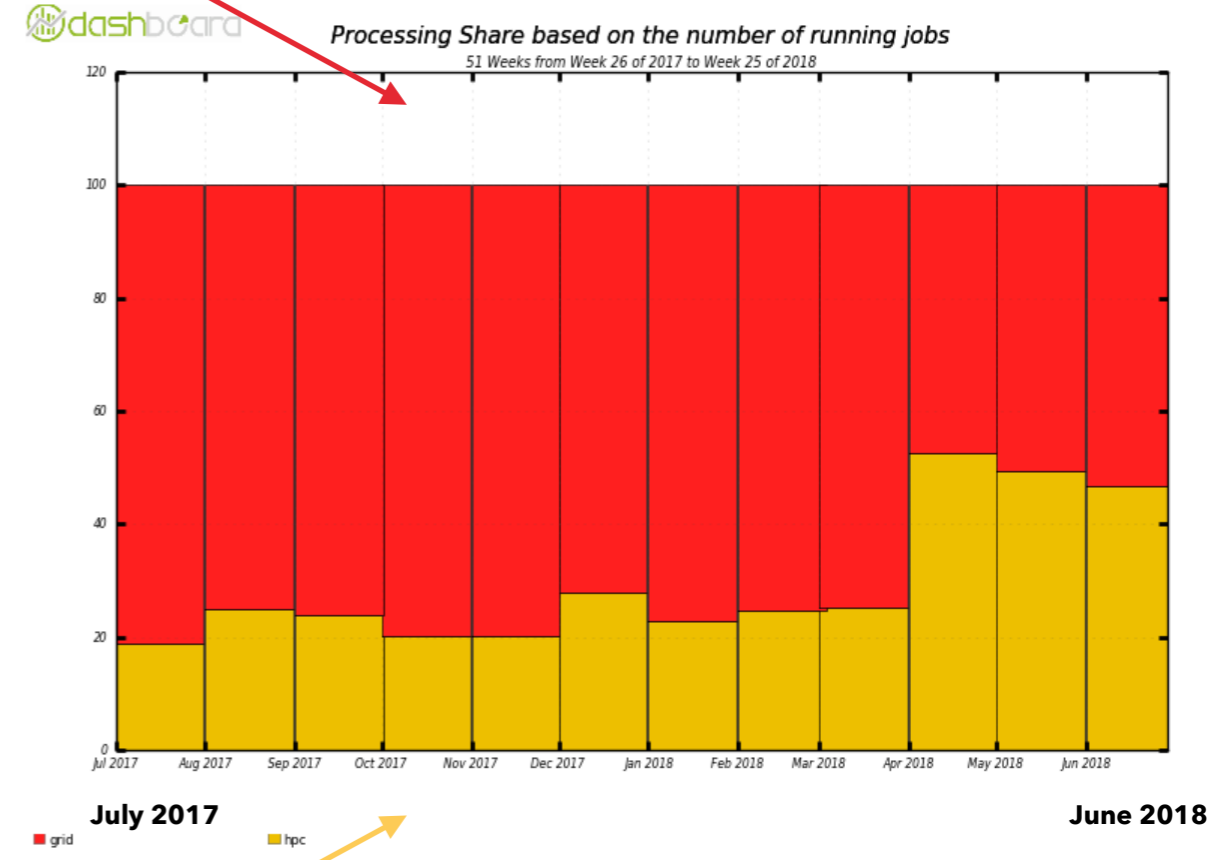
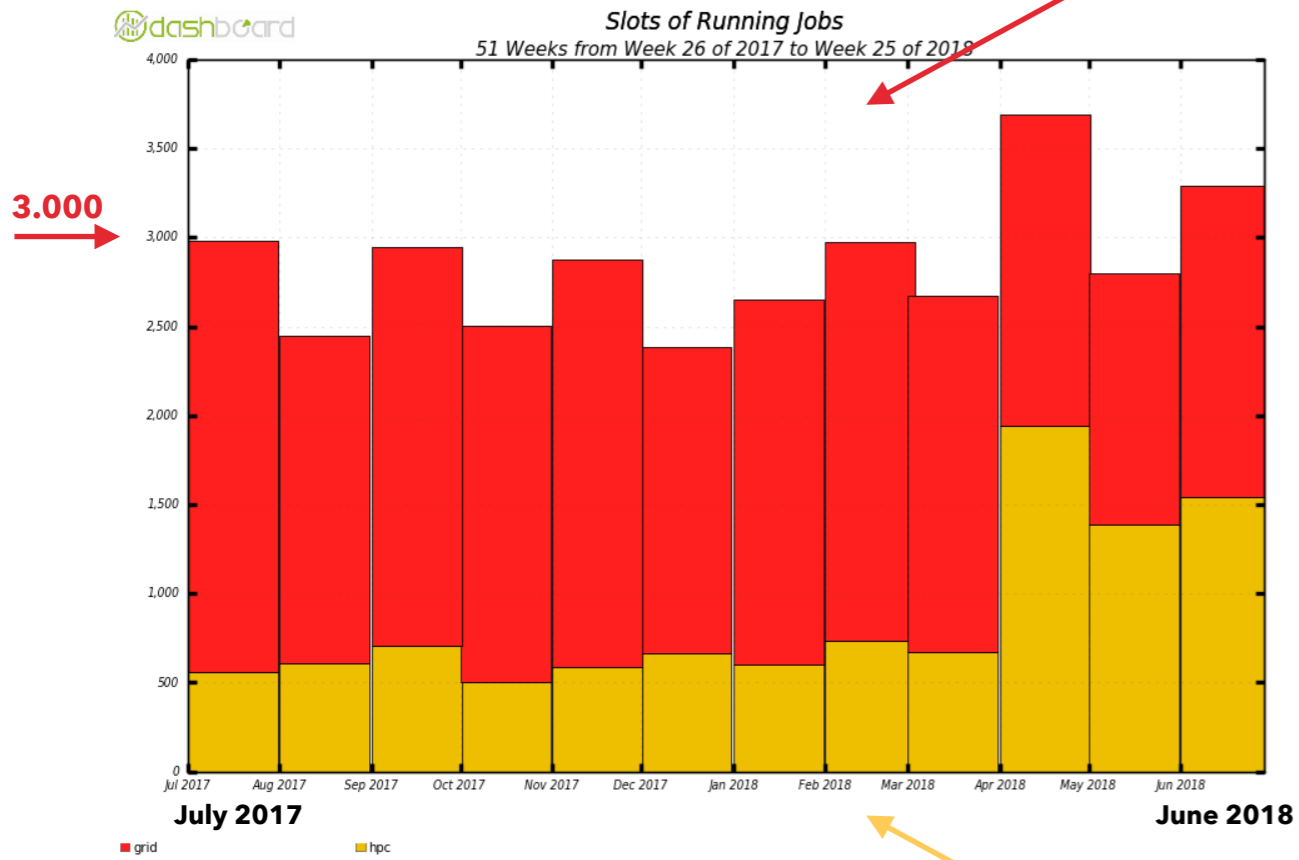


WLCG computing on HPC systems

## Used slots - 1 year

## Processing share - 1 year

x86\_64

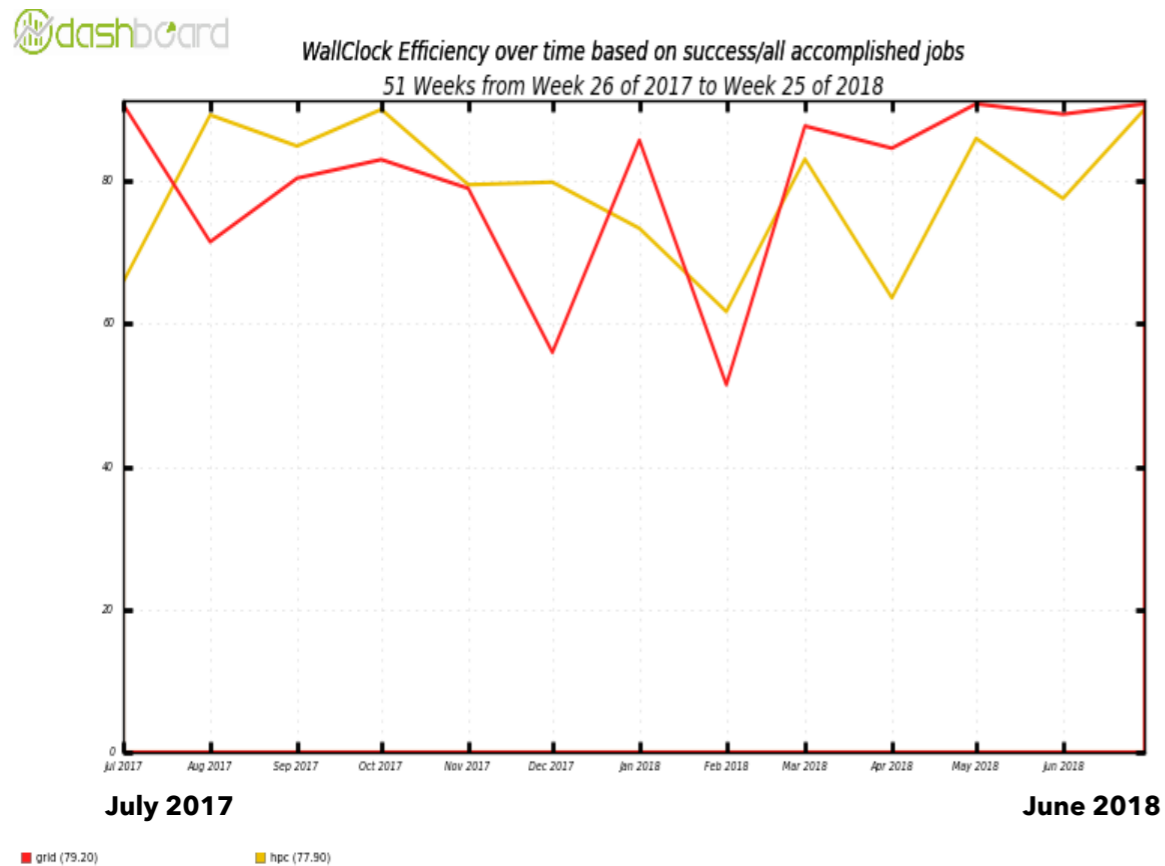


HPC

ATLAS only (~40% of the total)

WLCG computing on HPC systems

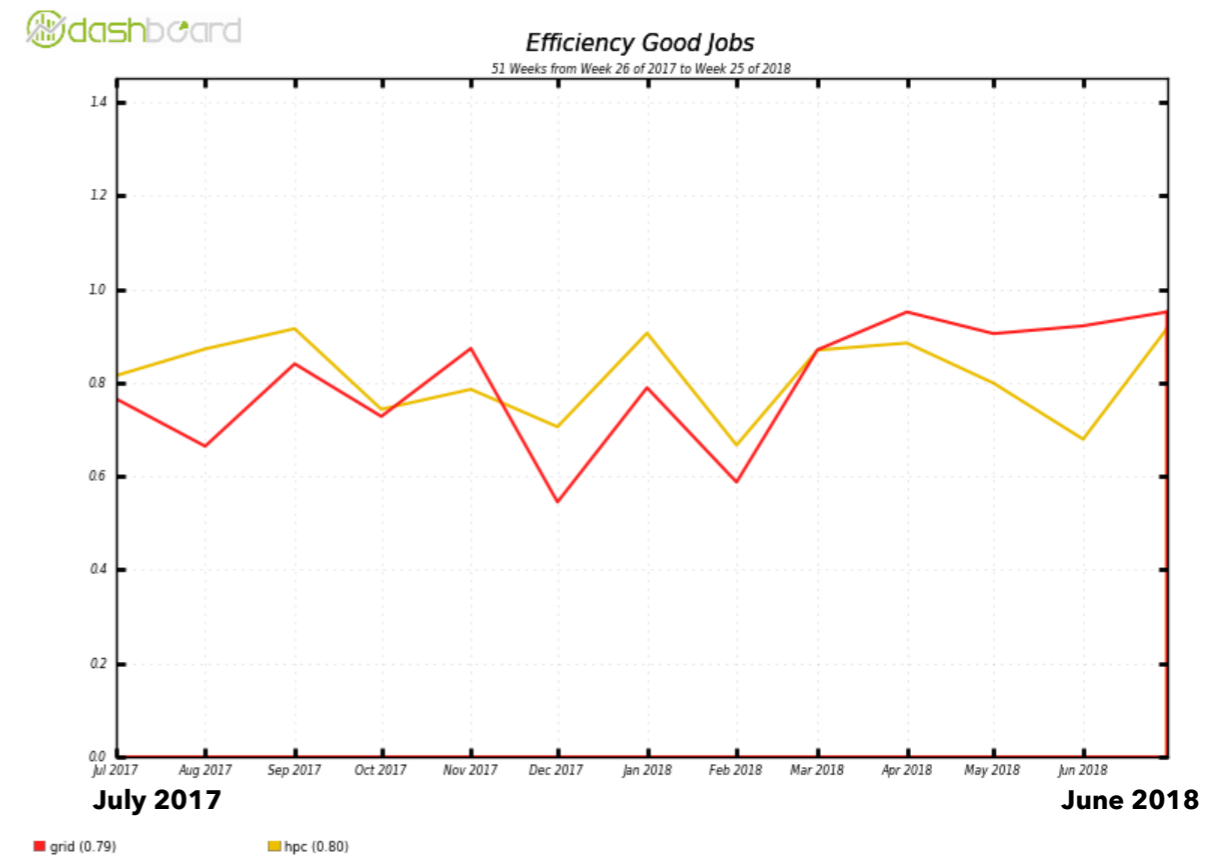
## WC efficiency (success/all)



**x86\_64: 79%**

**HPC: 78%**

## CPU/WC efficiency (good jobs)



**x86\_64: 79%**

**HPC: 80%**

**ATLAS only (~40% of the total)**

WLCG computing on HPC systems

- ▶ **Cost of resources marginally lower**
- ▶ **Comparable performance**
- ▶ **High integration costs and ongoing challenges**
- ▶ **Where is the big deal then?**

## WLCG computing on HPC systems

- ▶ **Cost of resources marginally lower**
- ▶ **Comparable performance**
- ▶ **High integration costs and ongoing challenges**
- ▶ **Where is the big deal then?**
  
- ▶ **Short term answer:**
  - There is NO big deal under such conditions and business model
  - No opportunistic usage means we get what we pay for (not elastic either)
  - It draws lots of attention, but the implementation effort and the operational pressure are considerably higher
  
- ▶ **HPC is sexy**
  - Could be made sexier

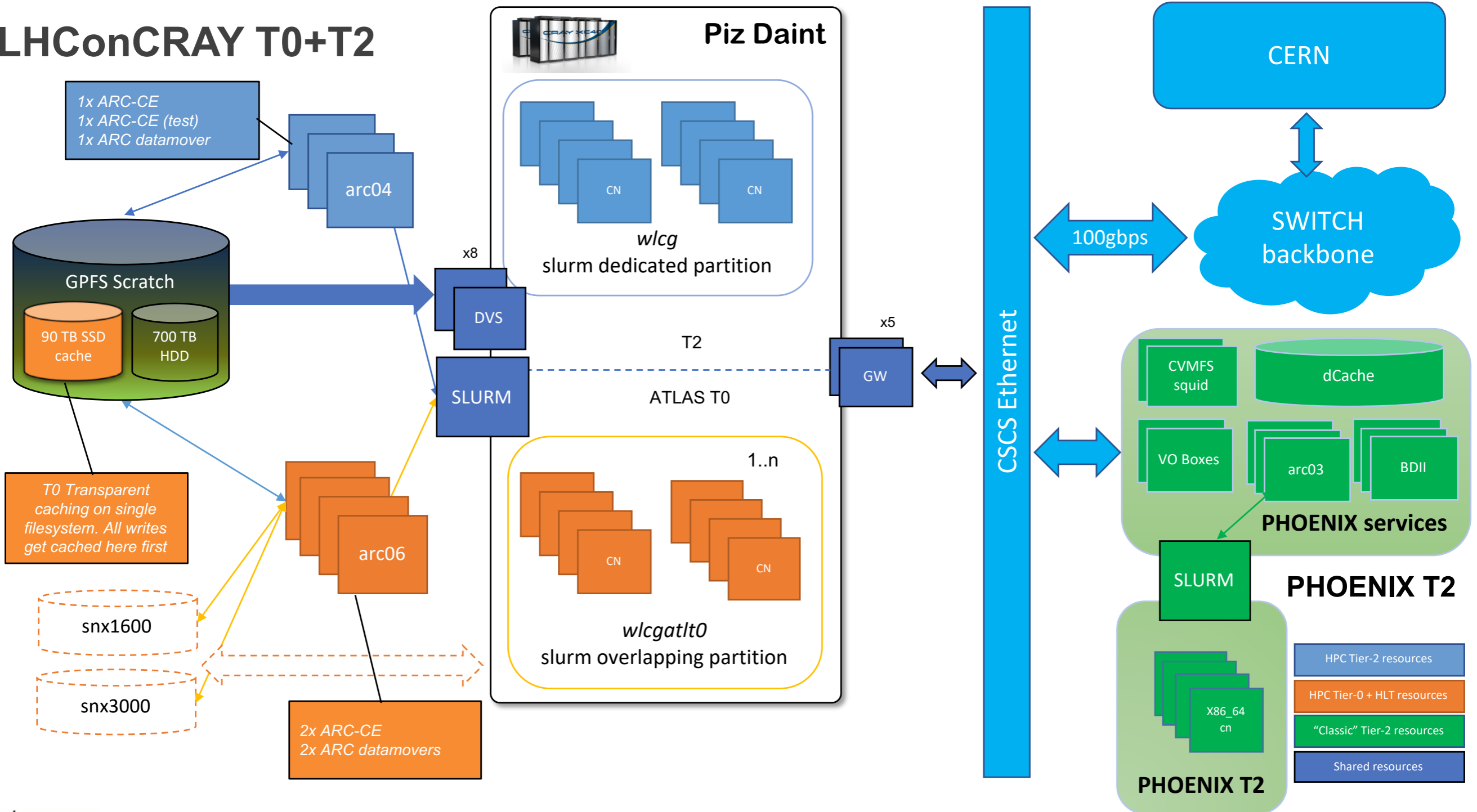
## WLCG computing on HPC systems

- ▶ **The architecture solutions may pave the way to future computing models**
  - E.G. shifter, CVMFS solutions
  - Debugging of DVS, potential in DWS
  - ...
- ▶ **And the other way around**
  - Community feedback to the next generation machine, on the way to exa-scale
- ▶ **The provisioning model is more flexible**
  - Still largely left to the good-will of the resource provider
  - Some freed up effort can be re-directed toward expert experiment support
- ▶ **The integration effort could pay back in the long term**
  - Opens doors to "*beyond pledge activities*"
  - E.G. elastic provisioning of resources for specific activities
  - One example: run ATLAS Tier-0 processing tasks (T0 spill-over or extension)
  - Several r&d projects are being proposed by ATLAS & CERN
  - Novel architectures, data lakes, code optimisation, porting to GPUs, etc.
  - Aiming at developing novel computing models

# T0 / T2 SHARED ARCHITECTURE

From T2 shared architecture to T2 / T0 shared

## LHCOnCRAY T0+T2



## Commissioning status

- ▶ **Dealt/dealing with teething problems**
  - Shortage of node memory for the type of workload
  - High and sustained I/O pressure
  - I/O patterns not usual for HPC environments
  - Exposing bugs in DVS, DWS Cray technologies
  - Adapt the workload to the architecture
  - ...
  
- ▶ **This is a step further to the HPC commissioning for Tier-2**
  - The Tier-0 workload is more demanding and resource hungry
  - On the Tier-2, the mix of workloads smooths out the edge requirements of a single one, both in terms of memory and I/O
  
- ▶ **Made good progress so far**
  - We are close to to be ready to go into production

## WLCG computing on HPC systems

- ▶ **LHC future computing needs should benefit from HPCs**
- ▶ **Integration with the experiment data processing frameworks is not trivial**
- ▶ **Efforts ongoing between the Swiss HEP community and CSCS since a few years, successfully integrated Piz Daint with the WLCG frameworks**
- ▶ **CERN and CSCS have mutual interest in developing new models**
- ▶ **R&D projects have been launched, ATLAS is heavily involved**
  - Novel architectures and software optimisation
- ▶ **Next generation HPCs should accommodate the requirements of the HEP community applications**



**Thank you for your attention!**

