
SWT2 Site Report

HEPiX
11/4/24

Zachary Booth on behalf of the SWT2 Collaboration

SWT2 Overview

- Southwest Tier-2 (SWT2) is a collaboration between the University of Oklahoma (OU) and the University of Texas at Arlington (UTA)
- Each site operates a data center on its campus
 - UTA: dedicated facility in the Chemistry & Physics Building
 - OU: within the OU Supercomputing Center for Education & Research (OSCER)
- Personnel: Zach Booth, Kaushik De (PI), Horst Severini, Mark Sosebee, Armen Vartapetian (K8s), Andrey Zarochentsev, Chris Walker
- Goal: **Provide robust & stable CPU cycles, data storage for ATLAS**



SWT2 Overview (II)

- The sites provide a mix of compute resources:
 - Traditional batch queues (slurm)
 - Kubernetes (K8s) cluster
 - Cloud resources (Google, since January '24)
- Total job slots / threads (both sites combined):
 - Baseline: ~25k (24k batch, 1k K8s)
 - Burst up: depending on workloads running in the Google PanDA queues, available / opportunistic job slots in OSCER
 - Routinely ~13k, with the potential ceiling much higher

SWT2 Overview (III)

- ATLAS utilizes the PanDA workload management system (WMS)
- PanDA “queues” at the sites are designed to accommodate the range of workloads PanDA distributes to the sites for execution:
 - OU_OSCER_ATLAS (primary queue for data production)
 - OU_OSCER_ATLAS_OPP (opportunistic access to OSCER slots)
 - SWT2_CPB (primary queue for both production and user analysis)
 - SWT2_CPB_K8S (Kubernetes sub-cluster)
 - SWT2_GOOGLE_ARM (ARM processors)
 - SWT2_GOOGLE_VHIMEM (workloads requesting higher amounts of RAM)
 - SWT2_GOOGLE_ULTRA (workloads requesting even greater amounts of RAM)
 - Queues for various testing purposes

Storage Overview (disk - no tape)

- SWT2 provides a total of ~13.7 PB of disk space
 - 13 PB at SWT2_CPB
 - 38 storage servers - Dell MD3460, R740xd2
 - 0.7 PB at OU_OSCER
 - 7 storage servers - Dell T630
- SWT2_CPB system is based on XRootD
 - Underlying filesystem: xfs
- OU_OSCER: currently XRootD, but will soon migrate to CephFS
- Google PanDA queues currently utilize the storage element (SE) at SWT2_CPB via the WAN
 - Work is ongoing to setup an SE within Google - should be operational soon
- An additional 5.8 PB will be added at SWT2_CPB this fall
 - Retire the oldest hardware
 - Net increase of ~3 PB

Networking Overview

- SWT2_CPB LAN is based on:
 - Two Dell S5232 core switches
 - Dell S4128 top-of-rack switches
- LAN was re-freshed in 2022

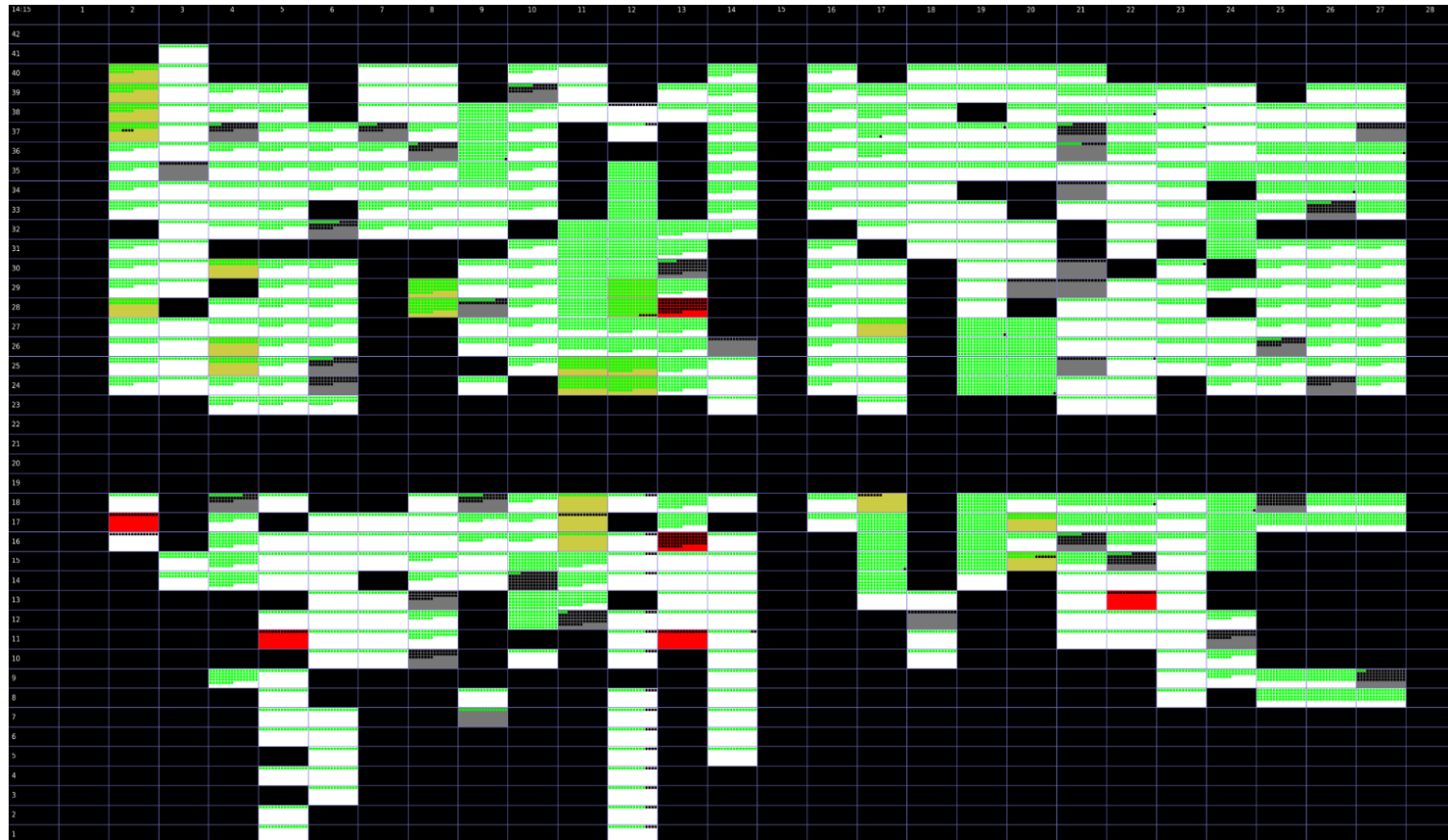
- WAN connectivity:
 - SWT2_CPB: 30 Gb/s dedicated maximum via LEARN => LHCONE
 - OU_OSCER: 100 Gb/s via OneNet => LHCONE
- External data transfers into / out of the SWT2 clusters:
 - four XRootD proxy data transfer nodes (DTN's) (SWT2_CPB)
 - one XRootD door (OU_OSCER)

Additional Services at the Sites

Various services operate to support data processing and storage activities:

- Compute Elements (CE's) via HTCondorCE
- XRootD redirectors - LAN access to the storage
- Slurm servers
- Squid caching proxies
- Nagios, custom dashboards for cluster monitoring

Map of Active Job Slots at SWT2_CPB



Evolution / Planned Improvements

- Maintain excellent performance, stability, scalability & diversity
 - XRootD, slurm, K8s, cloud (ARM, GPU, specialized PanDA queues, etc.)
- Optimize performance through technology evaluation
 - For example, batch vs K8s vs cloud vs ...
 - Alternate storage solutions like XRootD, EOS, Ceph
- Evaluate the correct mix of compute vs storage for Tier-2 evolution
 - WLCG pledges, ATLAS needs for the future
 - This then drives procurement decisions

Evolution / Planned Improvements (II)

- Migrate to Ceph-based storage at OU_OSCER
- Simplification of deployment, operations and monitoring
- Increase WAN bandwidth at SWT2_CPB for HL-LHC
 - Current bottleneck is a 30 Gb/s fiber link from the data center to the campus edge router
 - Discussions with campus networking staff are underway
- Alma9 OS upgrade:
 - Essentially done at OU_OSCER
 - In progress at SWT2_CPB

Thanks!

QUESTIONS?