



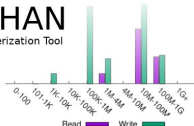
Integrating Darshan I/O Performance Monitoring Into SPOT

Author: Wesley Kwiecinski

Mentors: Rui Wang (ANL), Peter Van Gemmeren (ANL)

Advisor: Michael E. Papka (UIC / ANL)

DARSHAN
HPC I/O Characterization Tool



HEP-CCE/SOP



Overview

- Brief recap
- Summer work
- Integration into Software Performance Optimization Team (SPOT)
Performance Monitoring Board (PMB)
- Future Work

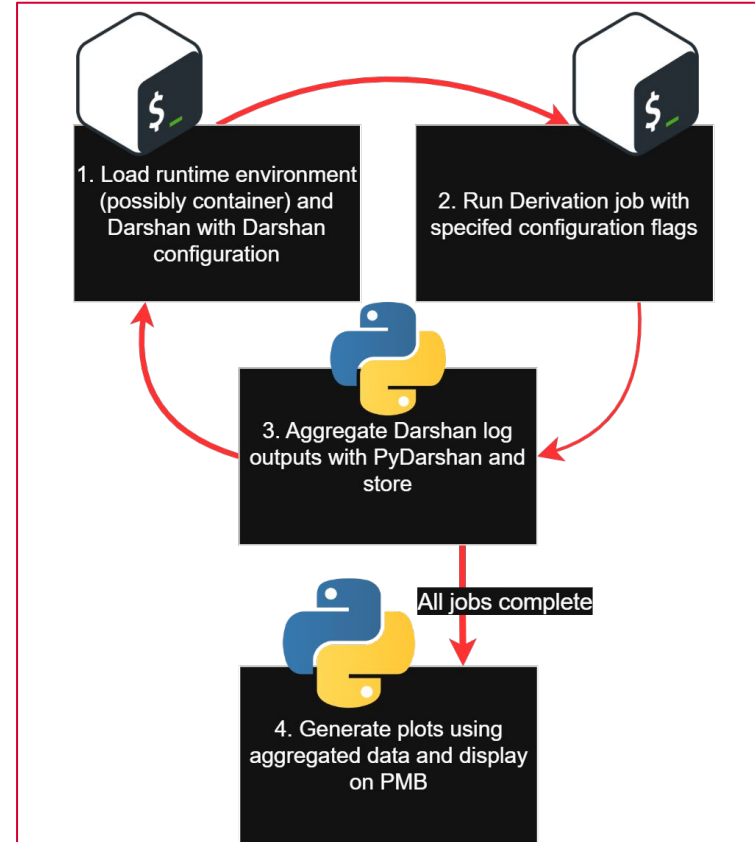
Brief Recap...

- Darshan would extend SPOT's tools with extra information for analysis
 - File timestamps, categories of I/O, modules for various I/O contexts (POSIX, STDIO, HDF5, etc)
 - Trace hotspots between Athena releases w/ timestamps (DXT)
 - Darshan extends SPOT with support for tracking individual workers
 - Darshan traces files per process
- Darshan has negligible overhead on HEP workflows^{*}

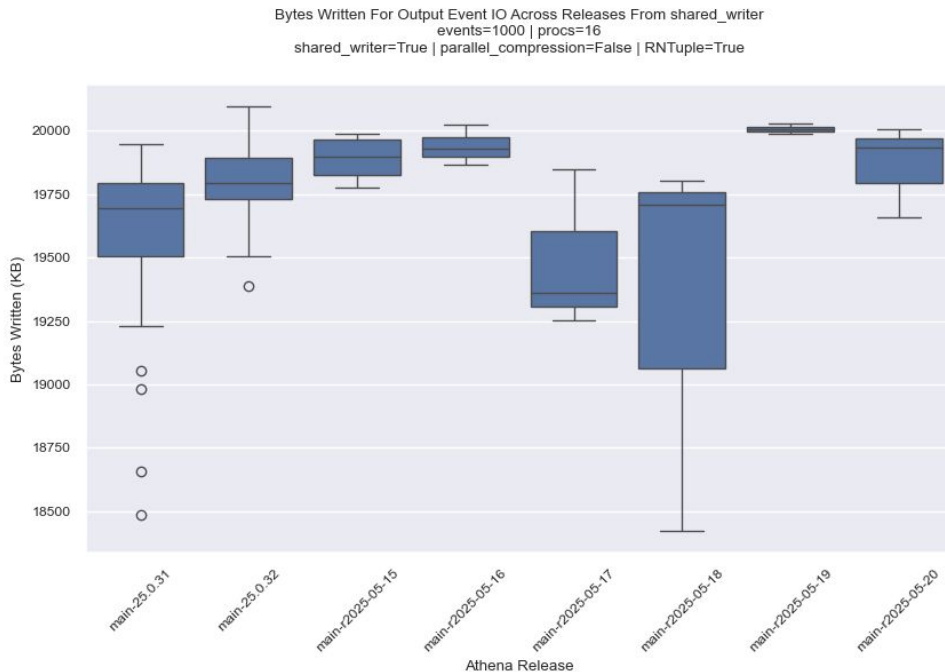
^{*} A. Vijayakumar. (2023). I/O Monitoring for portable HPC applications

Brief Recap...

- Use Cron + container + configurable bash script
 - Track I/O performance of Derivations w/ different configurations and releases



Brief Recap...



Difference in Average POSIX Reads & Writes for Event I/O files Per Process
events=1000 | procs=16
shared_writer=True | parallel_compression=False | RNTuple=True

Reads shared_writer

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (17)	0 (17)	0 (17)	0 (17)	0 (17)	0 (17)	0 (17)	0 (17)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Writes shared_writer

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
DAOD_PHYSLITE.pool.root.1	0 (11093)	2 (11095)	-4 (11091)	-5 (11086)	2 (11088)	6 (11094)	-4 (11090)	

Reads worker_4

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Writes worker_4

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Reads worker_9

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Writes worker_9

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Reads worker_14

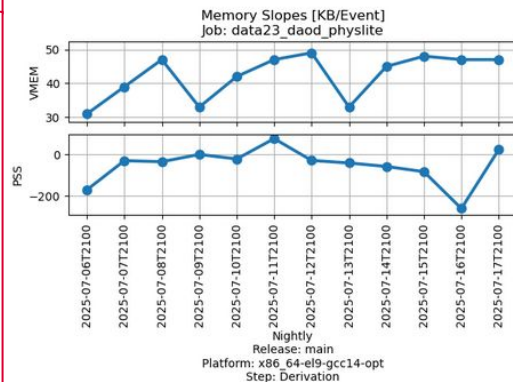
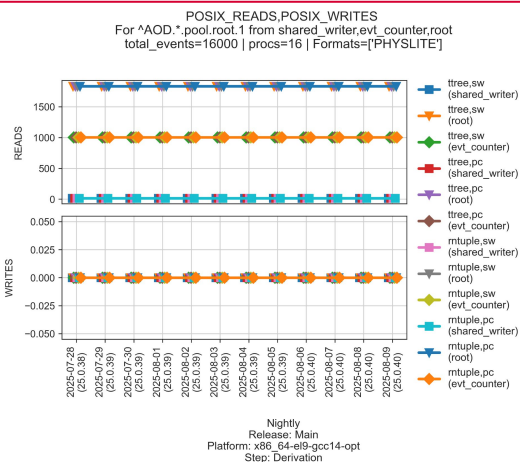
	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)	0 (38261)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Writes worker_14

	main-23.0.31	main-23.0.32	main-2023-05-15	main-2023-05-16	main-2023-05-17	main-2023-05-18	main-2023-05-19	main-2023-05-20
AOD_27162646_000001.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
DAOD_PHYSLITE.pool.root.1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

Athena Release

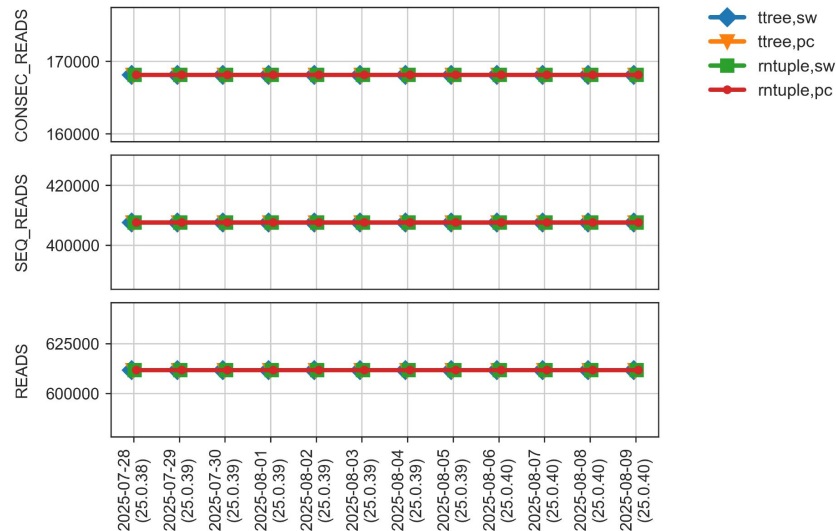
- Homogenization of plots to match SPOT PMB
- Internal work on plot generation scripts
- Gathering DXT data for more details on event I/O reads & writes
- Leaning into gathering data from 4 different Derivation configurations
- Integration with SPOT's PMB



New data plots

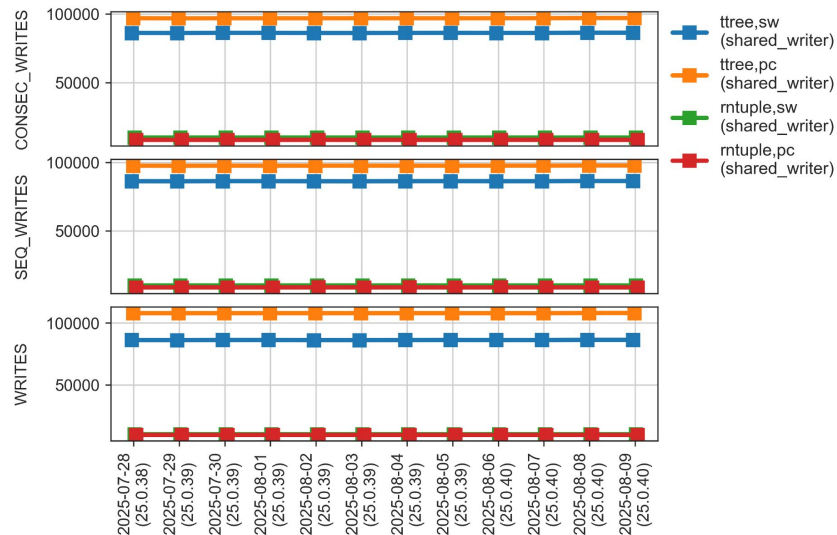
Summer Work

(many variables)
For ^AOD.*.pool.root.1 from all
total_events=16000 | procs=16 | Formats=['PHYSLITE']



Nightly
Release: Main
Platform: x86_64-el9-gcc14-opt
Step: Derivation

(many variables)
For ^DAOD_PHYSLITE.pool.root.1 from shared_writer
total_events=16000 | procs=16 | Formats=['PHYSLITE']



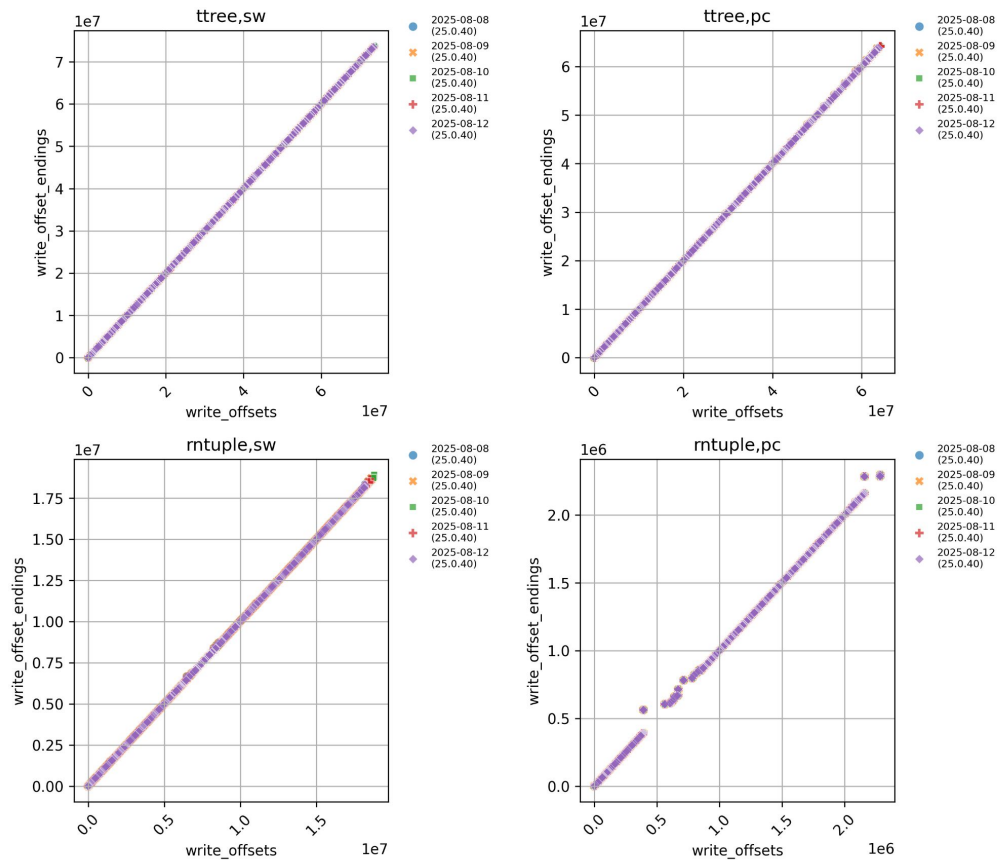
Nightly
Release: Main
Platform: x86_64-el9-gcc14-opt
Step: Derivation

Summer Work

DXT write_offsets vs write_offset_endings

file=.*DAOD_PHYSLITE.pool.root.1

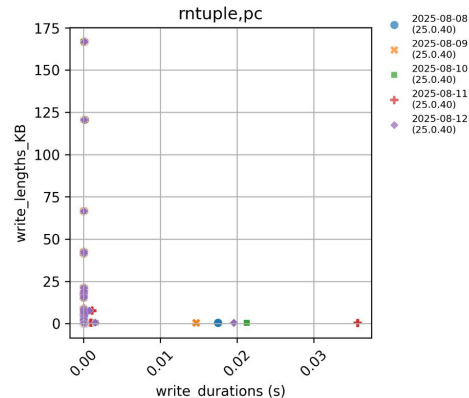
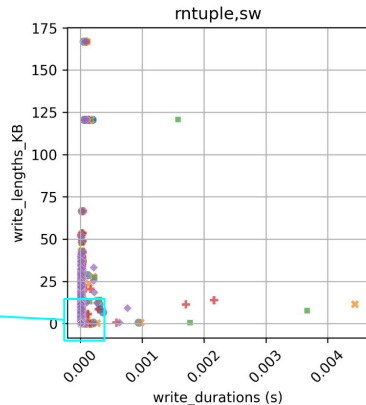
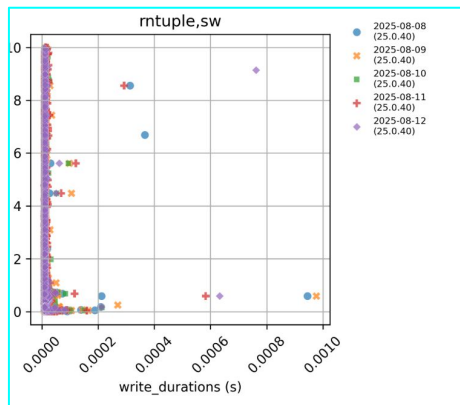
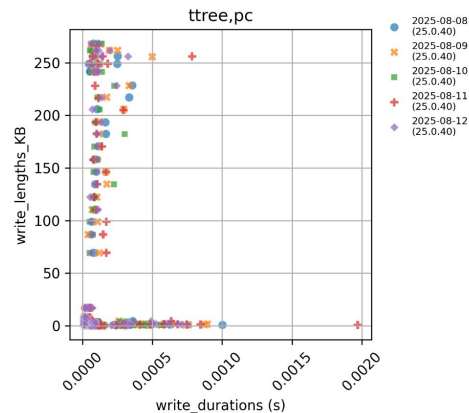
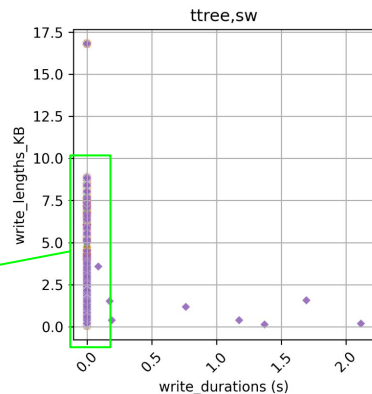
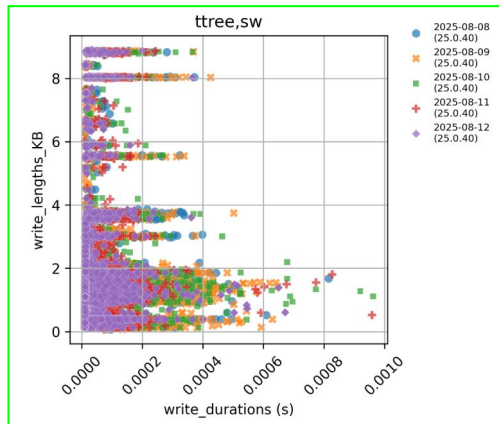
Total Events=16000



Release: main
Platform: x86_64-el9-gcc14-opt
Step: Derivation

Summer Work

DXT write_durations (s) vs write_lengths_KB
file=DAOD_PHYSLITE.pool.root.1
Total Events=16000



Release: main
Platform: x86_64-el9-gcc14-opt
Step: Derivation

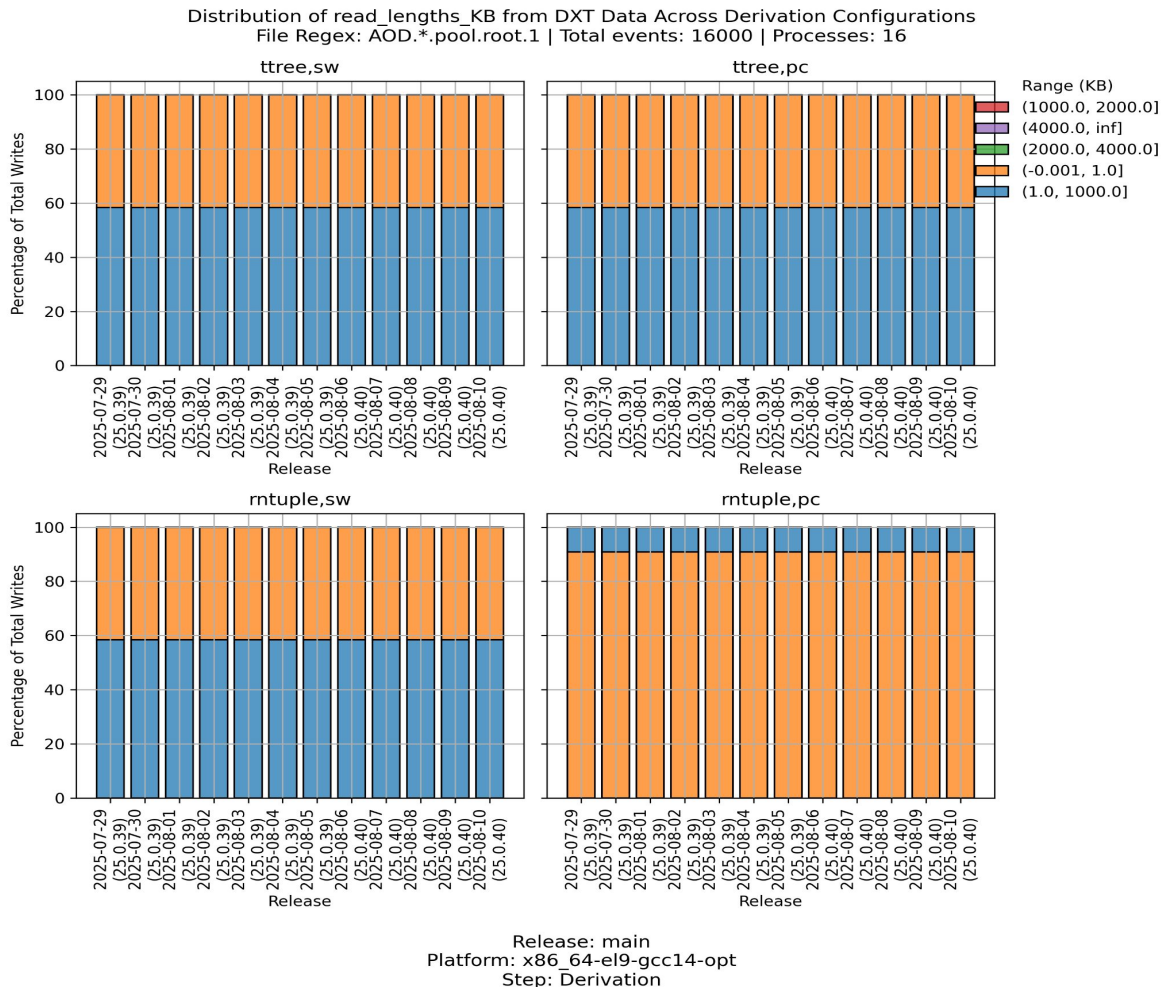
Summer Work

Distribution of write_lengths_KB from DXT Data Across Derivation Configurations
File Regexp: .*DAOD_PHYSLITE.pool.root.1 | Total events: 16000 | Processes: 16



Release: main
Platform: x86_64-el9-gcc14-opt
Step: Derivation

Summer Work



Integration of Scripts with SPOT PMB

- Spent time over the summer migrating job & plotting scripts to SPOT's PMB repository
- Merged with SPOTs PMB this week, starting to collect data
- Still need to include DXT plots in PMB, currently only using plots from posix module

Future Work

- Adjust DXT data collection
- Collect data from SPOT servers to identify possible I/O bottlenecks with Derivations
 - Still need to include DXT data, PMB will only show standard POSIX data for now
- Potential Master's Thesis topic
 - Currently exploring performance of small I/O in Derivations on machines with low vs high I/O contention
 - I/O Performance & I/O Contention on HPCs are generally well-studied, some case studies on different scientific software



Thank You!

Acknowledgements

- This work is a collaboration between the University of Illinois Chicago and Argonne National Laboratory under the C2-The-P2 fellowship
- This work uses resources at the Argonne Laboratory Computing Resource Center (LCRC)