

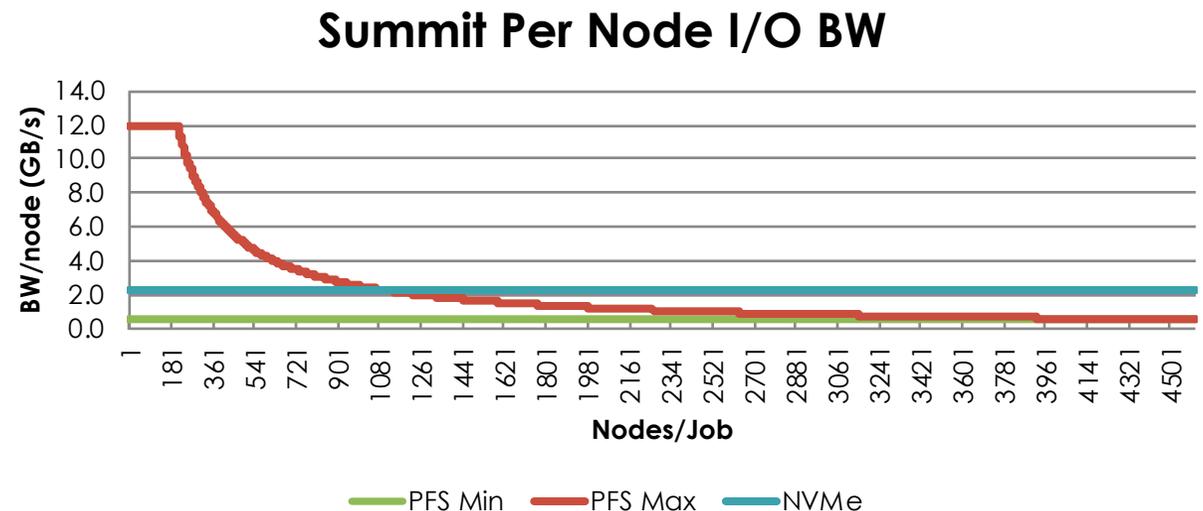
Summit Burst Buffer Libraries

Christopher Zimmer

ORNL is managed by UT-Battelle, LLC for the US Department of Energy

When to use the Burst Buffers

- Alpine GPFS Performance
 - Per node 12-14 GB/s
 - Aggregate 2.5 TB/s
 - Full system job will achieve 550 MB/s per node
- Node Local NVME
 - Samsung PM1725A
 - Write 2.1 GB/s
 - Read 5.5 GB/s
 - Scales linearly with Job Size



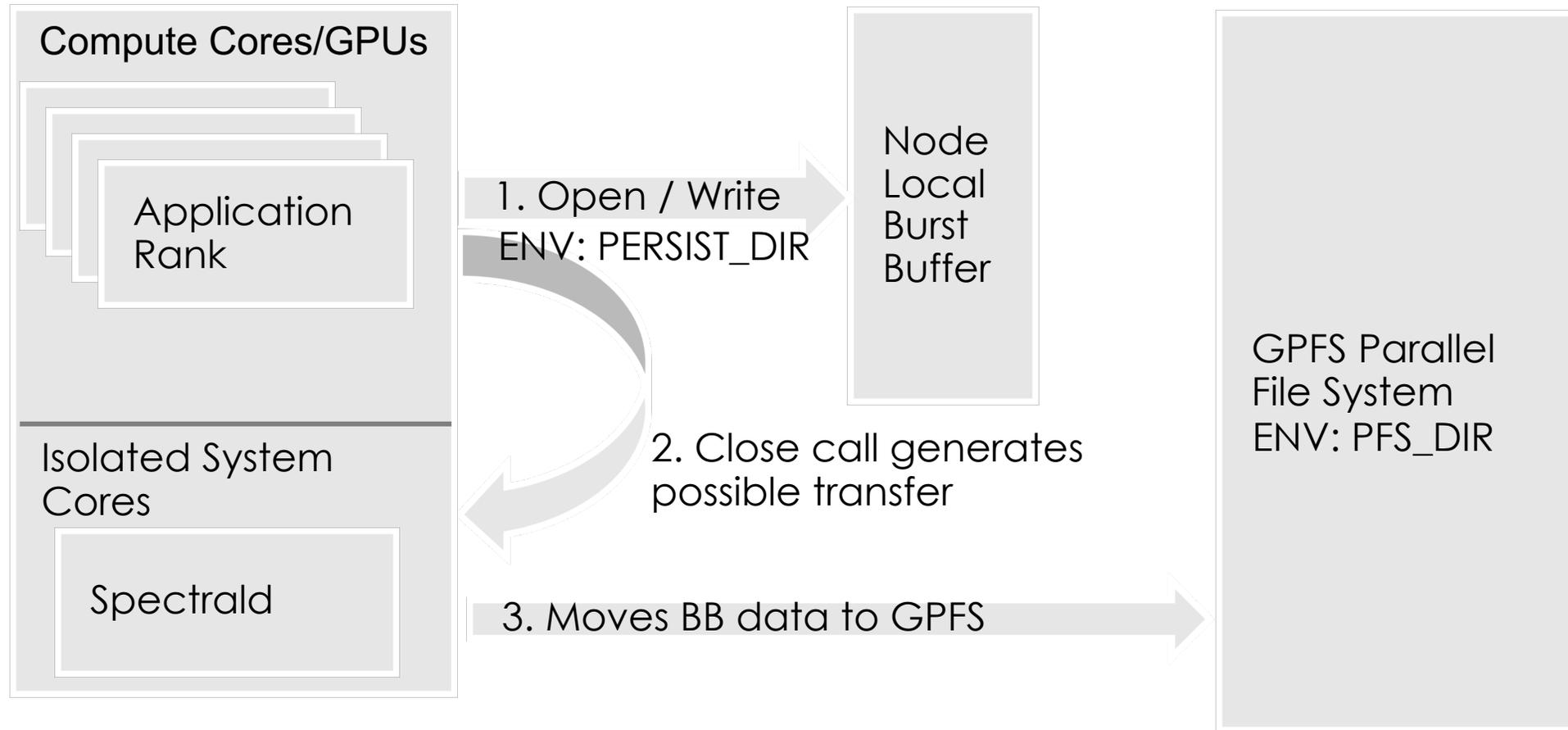
Modes of Use

- Scratch File System
 - Good for ML/DL Workloads
 - Self managed checkpointing
 - Ephemeral storage
- Spectral
 - File per process checkpoints
 - Iterative output
- SymphonyFS
 - Shared file output

Spectral

- On node copy agent
 - Runs on isolated cores as system agent
- Application Interface Transparent
 - Node code modifications (LD_PRELOAD)
 - Changes limited job scripts
 - Application only reasons about a single namespace
- Preserves portability with single namespace
- Non-shared files

Spectral Data Flow



Lammps Example

```
#Changes for Spectral
export PERSIST_DIR=${BBPATH}
export PFS_DIR=$PWD/restart
export OMPI_LD_PRELOAD_PREPEND=/ccs/techint/home/czimmer/Spectral-Test/build/src/libspectral.so

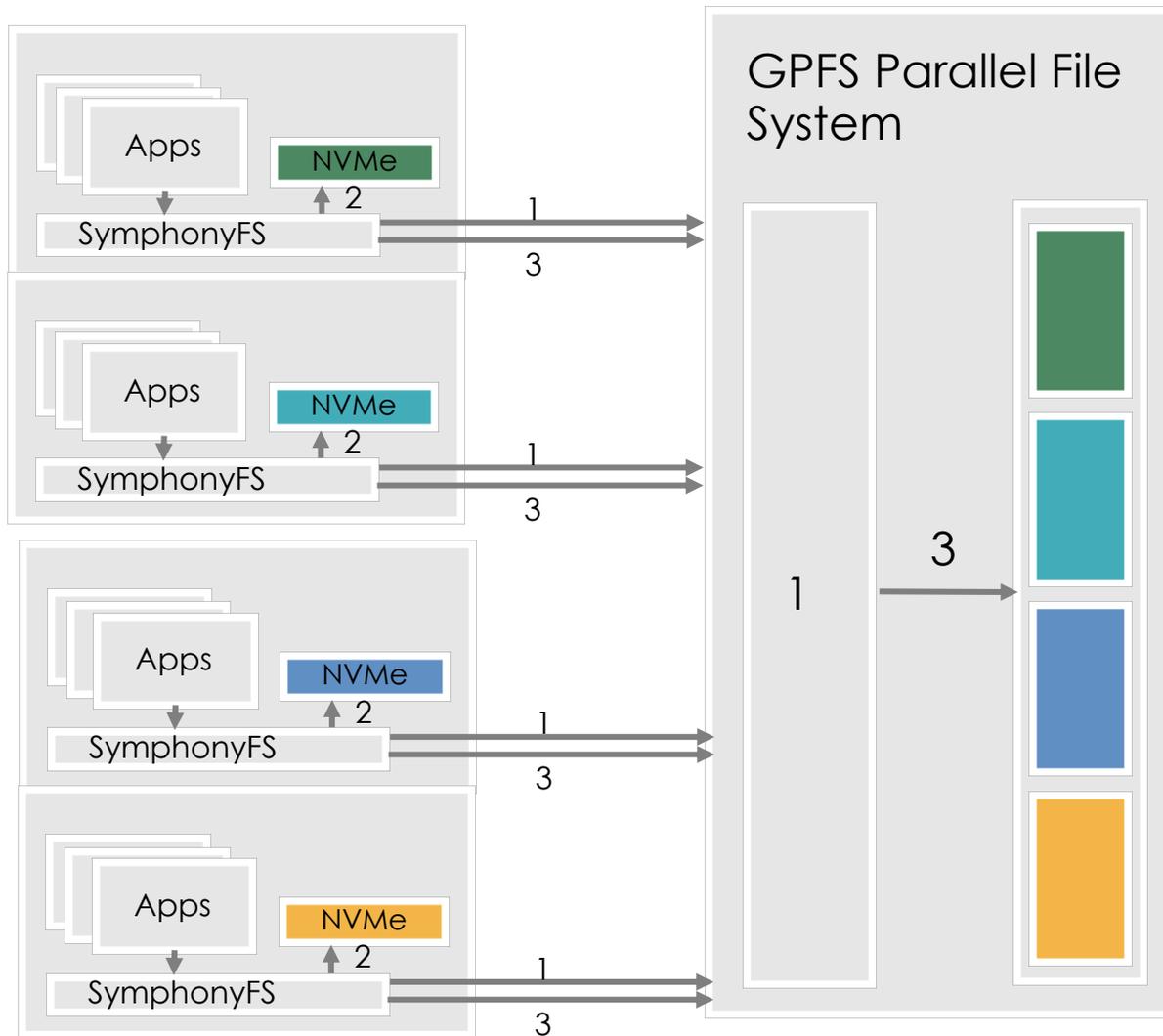
jsrun --nrs 16 --rs_per_host 1 --tasks_per_rs 6 ${BINARY} -v number_of_atoms ${natoms} -v output_dir ${PERSIST_DIR} < ${LAMMPS_INPUT_FILE}
```

- PERSIST_DIR sets where lammps actually writes to via config file
- When spectrald detects a write into PERSIST_DIR, the data is moved to PFS_DIR outside of the application
- OMPI_LD_PRELOAD_PREPEND will be set by module load

SymphonyFS

- Single (FUSE) name space
- File per process, shared file
- Operational model:
 - Uses NVMe as a write-back extent cache
 - Log based file system
 - Reconstructs file on parallel file system

SymphonyFS Data Flow



1. File open (Meta data handled through GPFS)'
2. Apps write data (buffered into NVMe's)
3. Upon file close and flush. SymphonyFS reconstitutes file on GPFS.

SymphonyFS

- Limitations
 - Non-overlapping writes
 - Read after write
 - Must be flushed to parallel file system
 - Reads come from parallel file system

Questions/Interest in early access?

- Spectral - zimmercj@ornl.gov
- SymphonyFS – brumgardcd@ornl.gov