

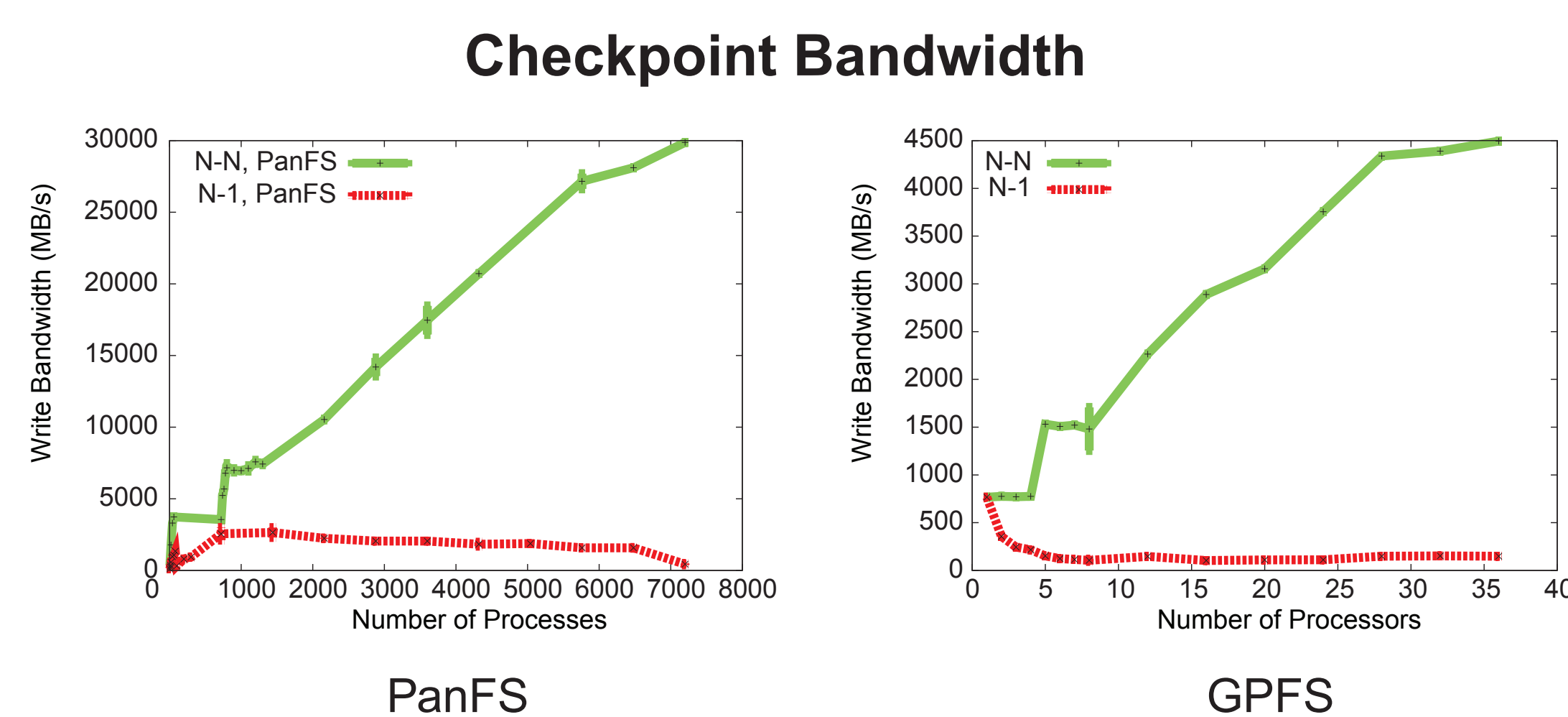
PLFS: A Checkpoint Filesystem for Parallel Applications

John Bent*, Garth Gibson†, Gary Grider*, Ben McClelland*, Paul Nowoczynski‡, James Nunez*, Milo Polte†, Meghan Wingate*

*Los Alamos National Laboratory †Carnegie Mellon University ‡Pittsburgh Supercomputing

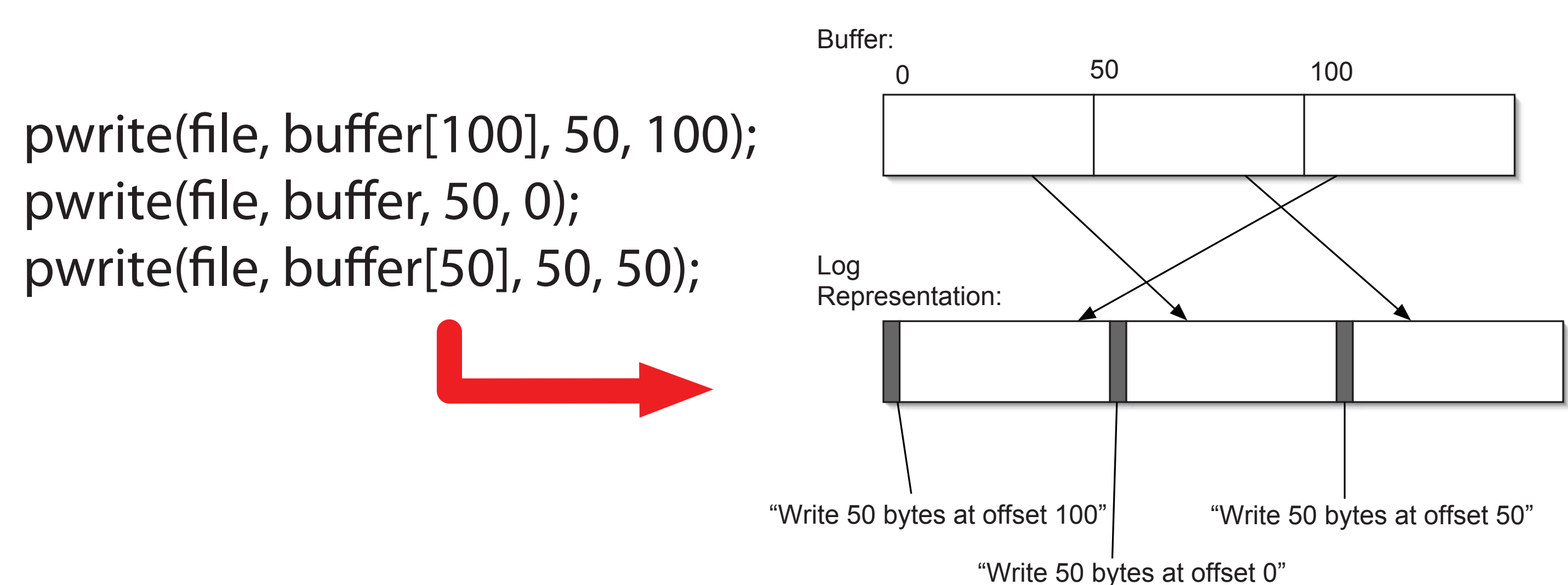
Problem

- Many important scientific applications create checkpoints using small, strided, concurrent writes to a shared file (N-1 checkpointing)
- Filesystems perform best on non-concurrent sequential workloads, such as N-N checkpointing
- Small-strided writes to a shared file often suffer from seeks and false sharing
- Unfortunately, we can't change the applications, but we can modify our filesystems



Previous Work

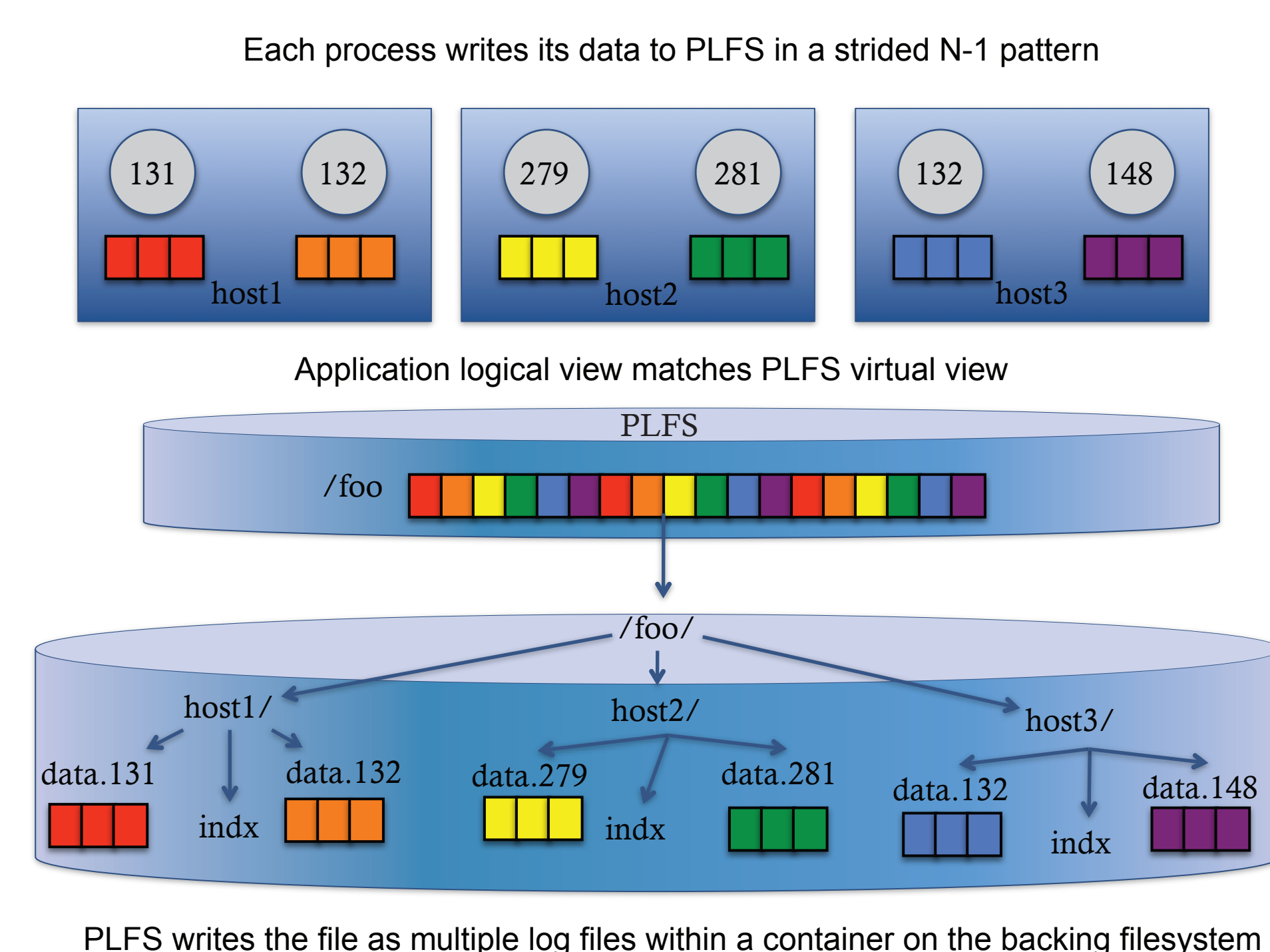
- "Log-structured Files for Fast Checkpointing"
- CMU students modified a parallel filesystem (PVFS2) to write shared files in a log representation
 - Required server modification
 - Only works with one filesystem
 - All writes are appended so no seeks, but clients still concurrently access a single file



PLFS – Parallel Log-Structured File System

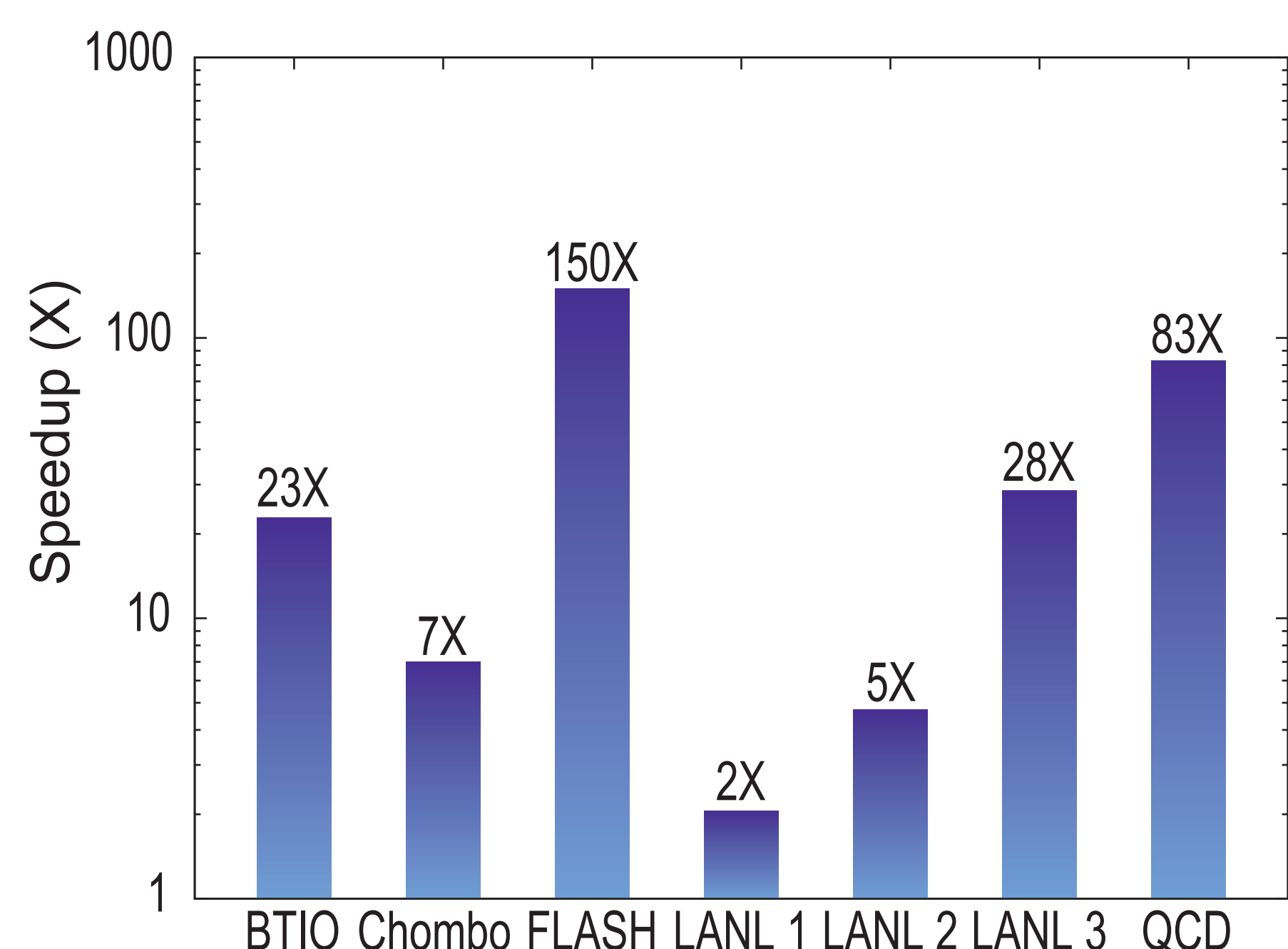
- A project developing a filesystem level improvement to N-1 checkpointing, led by John Bent (LANL)
- FUSE based filesystem mounted on top of any existing parallel filesystem on clients
- Decouples a concurrent N-1 checkpoint into a non-concurrent N-N checkpoint
- Redirects strided writes from multiple processes accessing a single file to sequential writes to data logs and index files

Layout of a PLFS Container



PLFS Speedup

- 2x – 150x speedups for important HPC applications at LANL scale!



with information from LANL Technical Release LA-UR 09-02117

Future Work

- Metadata optimizations
- Read path optimizations
 - Currently can perform poorly if read and write access patterns differ
 - Metadata servers could keep in-memory maps of the indices
- Specialization for patterns other than N-1
- PLFS-like approaches to distribute directories on filesystems that store each on a single metadata server