

Accelerating your application I/O with UnifyFS

HPC-IODC 2020 June 25, 2020







LLNL: KATHRYN MOHROR (PI), ADAM MOODY, CAMERON STANAVIGE, TONY HUTTER

ORNL: SARP ORAL (ORNL-PI), HYOGI SIM, FEIYI WANG, MIKE BRIM, SWEN BOEHM

NCSA: CELSO MENDES, CRAIG STEFFEN



What is UnifyFS?

- Simply put, it's a file system for burst buffers
- Our goal is to make using burst buffers on exascale systems as *easy* as writing to the parallel file system and orders of magnitude *faster*

.nt main(int argc, char **argv) { MPI_Init(argc, argv);	<pre>void checkpoint(void) { int rank;</pre>
for (t = 0; t < TIMESTEPS; t++) {	<pre>MPI_Comm_rank(MPI_COMM_WORLD, &rank);</pre>
/* do work */	-// file = "/pfs/shared.chpt"; file = "/unifyfs/shared.ckpt";
}	<pre>File *fs = fopen(file, "w");</pre>
<pre>MPI_Finalize(); return 0; The only required</pre>	if (rank == 0) fwrite(header,, fs);
change is to use /unifyfs instead of /pfs	<pre>long offset = header_size +</pre>



What is UnifyFS?

- Simply put, it's a file system for burst buffers
- Our goal is to make using burst buffers on exascale systems as *easy* as writing to the parallel file system and orders of magnitude *faster*
- Results on Summit show scalable write performance for UnifyFS with shared files on burst buffers

UnifyFS Write Performance on Summit BBs



Writing data to the parallel file system is expensive



HPC Storage is becoming more complex



HPC Storage is becoming more complex



HPC Storage is becoming more complex



LLNL-PRES-811832

https://github.com/LLNL/UnifyFS





**Measurements of local storage performance taken with SCR (The Scalable Checkpoint/Restart Library) https://github.com/llnl/scr

https://github.com/LLNL/UnifyFS



• Sharing files on node-local burst buffers is not natively supported



- UnifyFS makes sharing files easy
- UnifyFS presents a shared namespace across distributed storage
- Used directly by applications or indirectly via higher level libraries like VeloC, MPI-IO, HDF5, PnetCDF, ADIOS, etc.
- UnifyFS is fast
- Tailored for specific HPC workloads, e.g., checkpoint/restart, visualization output
- Each UnifyFS instance exists only within a single job, no contention with other jobs on the system

UnifyFS is designed to work completely in user space for a single user's job





```
int main(int argc, char* argv[]) {
  MPI_Init(argc, argv);
  read_input();
  for(int t = 0; t < TIMESTEPS; t++)</pre>
    /* ... Do work ... */
    /* ... Synchronization ... */
      write_output();
  }
  MPI_Finalize();
  return 0;
```





```
int main(int argc, char* argv[]) {
  MPI_Init(argc, argv);
  read_input();
  for(int t = 0; t < TIMESTEPS; t++)</pre>
    /* ... Do work ... */
    /* ... Synchronization ... */
      write_output();
  }
  MPI_Finalize();
  return 0;
```



• Reads and writes occur in distinct phases



```
int main(int argc, char* argv[]) {
 MPI_Init(argc, argv);
  read_input();
  for(int t = 0; t < TIMESTEPS; t++)</pre>
    /* ... Do work ... */
    /* ... Synchronization ... */
      write_output();
  }
  MPI Finalize();
  return 0;
```



- Reads and writes occur in distinct phases
- Reads and writes are performed to regular offsets in files (not random writes or reads)



```
int main(int argc, char* argv[]) {
  MPI_Init(argc, argv);
  read_input();
  for(int t = 0; t < TIMESTEPS; t++)</pre>
    /* ... Do work ... */
    /* ... Synchronization ... */
      write output();
  }
  MPI Finalize();
  return 0;
```



- Reads and writes occur in distinct phases
- Reads and writes are performed to regular offsets in files (not random writes or reads)
- When reading, a process will most likely access data it wrote



```
int main(int argc, char* argv[]) {
  MPI_Init(argc, argv);
  read_input();
  for(int t = 0; t < TIMESTEPS; t++)</pre>
    /* ... Do work ... */
    /* ... Synchronization ... */
      write output();
  }
  MPI Finalize();
  return 0;
```

- Reads and writes occur in distinct phases
- Reads and writes are performed to regular offsets in files (not random writes or reads)
- When reading, a process will most likely access data it wrote
- These behaviors are typical of common HPC workloads
 - Checkpoint/restart, periodic output, producer/consumer

https://github.com/LLNL/UnifyFS

UnifyFS files are globally visible after "lamination"

- Based on typical I/O behavior in HPC applications we utilize "lamination semantics"
- Before lamination processes may not see updates written by processes on another node
- Once processes are done modifying a file, they initiate lamination
 - The lamination process renders your file read-only and synchronizes file data across nodes in your job
 - Now any process on any node can read the final state of the file
 - And you can transfer data to external storage



UnifyFS files are globally visible after "lamination"

- Based on typical I/O behavior in HPC applications we utilize "lamination semantics"
- Before lamination processes may not see updates written by processes on another node
- Once processes are done modifying a file, they initiate lamination
 - The lamination process renders your file read-only and synchronizes file data across nodes in your job
 - Now any process on any node can read the final state of the file
 - And you can transfer data to external storage





- get and build UnifyFS
- build your application to use UnifyFS
- run your application with UnifyFS
- move your data between UnifyFS and the parallel file system





Got Spack?

\$ spack install unifyfs

\$ spack load unifyfs

https://github.com/spack/spack

How do I modify my MPI application for UnifyFS?

```
int main(int argc, char * argv[]) {
   FILE *fp;
   //program initialization
   //MPI setup
   //perform I/O
   fp = fopen("/lustre/out.txt", "w");
   fprintf(fp, "Hello World! I'm rank %d", rank);
   fclose(fp);
```

//clean up
return 0;

How do I modify my MPI application for UnifyFS?

```
int main(int argc, char * argv[]) {
    FILE *fp;
    //program initialization
    //MPI setup
```

```
//perform I/O
fp = fopen("/unifyfs/out.txt", "w");
fprintf(fp, "Hello World! I'm rank %d", rank);
fclose(fp);
```

```
// "laminate" the file to indicate to UnifyFS you
are done modifying this file
    chmod("/unifyfs/out.txt", 0444);
```

//clean up
return 0;

How do I modify my MPI application for UnifyFS?

int main(int argc, char * argv[]) {
 FILE *fp;
 //program initialization
 //MPI setup

//perform I/0
fp = fopen("/unifyfs/out.txt", "w");
fprintf(fp, "Hello World! I'm rank %d",
fclose(fp);

// "laminate" the file to indicate to Uni
are done modifying this file
 chmod("/unifyfs/out.txt", 0444);

//clean up
return 0;

"chmod" is currently supported method for lamination

- (Near-)future methods we plan to support
 - Laminate on close()
 - Laminate on unmount() (either explicit API call or auto-mount)
 - Laminate with API call, e.g., unifyfs_laminate
- Expect to be able to specify which method you want when you start up UnifyFS

unify_{FS}



How does UnifyFS intercept I/O calls?

• Static Linking

\$ mpicc -o hello `<unifyfs>/bin/unifyfs-config --pre-ld-flags` \
 hello.c `<unifyfs>/bin/unifyfs-config --post-ld-flags`

- Dynamic Linking (Recommended method)
 - Using Spack makes this very easy! (just "spack load unifyfs")
 - We use the Gotcha library for dynamic interception

\$ mpicc -o hello -L<unifyfs>/lib hello.c -lunifyfs_gotcha



• Easiest method: Start & stop UnifyFS in your batch script

allocate nodes and options for resource manager
#BSUB -nnodes 8 ...



- Easiest method: Start & stop UnifyFS in your batch script
 - Command 'unifyfs start' launches UnifyFS for your job and sets up the file system

allocate nodes and options for resource manager
#BSUB -nnodes 8 ...

shell command portion of batch script
unifyfs start --share-dir=/shared/file/system/path --mount=/unifyfs



- Easiest method: Start & stop UnifyFS in your batch script
 - Command 'unifyfs start' launches UnifyFS for your job and sets up the file system

allocate nodes and options for resource manager
#BSUB -nnodes 8 ...

shell command portion of batch script
unifyfs start --share-dir=/shared/file/system/path --mount=/unifyfs



- Easiest method: Start & stop UnifyFS in your batch script
 - Command 'unifyfs start' launches UnifyFS for your job and sets up the file system
 - Run your command as usual and use the path /unifyfs to direct data to UnifyFS

allocate nodes and options for resource manager
#BSUB -nnodes 8 ...

```
### shell command portion of batch script
unifyfs start --share-dir=/shared/file/system/path --mount=/unifyfs
jsrun -p 4096 ./hello
```



- Easiest method: Start & stop UnifyFS in your batch script
 - Command 'unifyfs start' launches UnifyFS for your job and sets up the file system
 - Run your command as usual and use the path /unifyfs to direct data to UnifyFS

allocate nodes and options for resource manager
#BSUB -nnodes 8 ...

```
### shell command portion of batch script
unifyfs start --share-dir=/shared/file/system/path --mount=/unifyfs
jsrun -p 4096 ./hello
```



- Easiest method: Start & stop UnifyFS in your batch script
 - Command 'unifyfs start' launches UnifyFS for your job and sets up the file system
 - Run your command as usual and use the path /unifyfs to direct data to UnifyFS
 - Command 'unifyfs terminate' cleans up the UnifyFS file system and tears it down

```
### allocate nodes and options for resource manager
#BSUB -nnodes 8 ...
```

```
### shell command portion of batch script
unifyfs start --share-dir=/shared/file/system/path --mount=/unifyfs
jsrun -p 4096 ./hello
unifyfs terminate
```



- Easiest method: Start & stop UnifyFS in your batch script
 - Command 'unifyfs start' launches UnifyFS for your job and sets up the file system
 - Run your command as usual and use the path /unifyfs to direct data to UnifyFS
 - Command 'unifyfs terminate' cleans up the UnifyFS file system and tears it down

```
### allocate nodes and options for resource manager
#BSUB -nnodes 8 ...
```

shell command portion of batch script
unifyfs start --share-dir=/shared/file/system/path --mount=/unifyfs
jsrun -p 4096 ./hello
unifyfs terminate



How do I move my data into and out of UnifyFS?

- Three ways to move data between UnifyFS and the parallel file system
- UnifyFS transfer program
 - jsrun -r1 transfer /unifyfs/file1 /scratch/mydir/file1
- UnifyFS transfer API
 - unifyfs_transfer_file_parallel("/unifyfs/out.txt", "/scratch/out.txt");
- Stage in and out options with UnifyFS command "unifyfs start"
 - unifyfs start --stage-in=/path/to/parallel/file/system/inputs --stageout=/path/to/parallel/file/system/outputs

unify_{FS}



- New features and improvements being added all the time
- Documentation and user support
 - User Guide: <u>http://unifyfs.readthedocs.io</u>
 - Down at the bottom, where it says "Read the Docs" chose "v: latest" to get the latest information
 - unifyfs@exascaleproject.org
- Example programs available in the "examples" directory in the source code

	 Docs » UnifyFS: A file system for burst buffers
nifyFS	UnifyFS: A file system for burst buffers
	User Guide
,t	Overview
	• High Level Design
	Definitions
	dol o
	Run or Job Step
	Assumptions
	 Application Behavior
	 Consistency Model
	 File System Behavior
	 System Characteristics
	Build & I/O Interception
	 UnifyFS Build Configuration Options
	 How to Build UnifyFS
	 I/O Interception
	Mounting UnifyFS
	 Mounting
	 Unmounting
	UnifyFS Configuration
	 o unifyfs.conf
	 Environment Variables
	 Command Line Options
X	Starting & Stopping in a Job
v: latest	- Starting Unit ES



- Our goal is to provide **easy, portable,** and **fast** support for burst buffers for ECP applications
- We need early users
- What features are most important to you
- Available on github: <u>https://github.com/LLNL/UnifyFS</u>
 - MIT license
- Documentation and user support
 - User Guide: <u>http://unifyfs.readthedocs.io</u>
 - unifyfs@exascaleproject.org









This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.



This research was supported by the Exascale Computing Project (17-SC-20-SC), a joint project of the U.S. Department of Energy's Office of Science and National Nuclear Security Administration, responsible for delivering a capable exascale ecosystem, including software, applications, and hardware technology, to support the nation's exascale computing imperative.

EXASCALE COMPUTING PROJECT