

Evaluating Integration and Performance of Containerized Climate Applications on a HPE Cray System

Subil Abraham

Oak Ridge National Laboratory
Oak Ridge, TN, USA
abrahams@ornl.gov

Ryan Prout

Oak Ridge National Laboratory
Oak Ridge, TN, USA
proutrc@ornl.gov

Thomas Robinson

*SAIC/NOAA/OAR/
Geophysical Fluid Dynamics Laboratory*
Princeton, NJ, USA
thomas.robinson@noaa.gov

Chris Blanton

NOAA/Geophysical Fluid Dynamics Laboratory
Princeton, NJ, USA
chris.blanton@noaa.gov

Luis Sal-bey

*SAIC/NOAA/OAR/
Geophysical Fluid Dynamics Laboratory*
Princeton, NJ, USA
luis.sal-bey@noaa.gov

Matthew Davis

Oak Ridge National Laboratory
Oak Ridge, TN, USA
davisjmj@ornl.gov

Abstract—Containers have taken over large swaths of cloud computing as the most convenient way of packaging and deploying applications. The features that containers offer for packaging and deploying applications translate to High Performance Computing (HPC) as well. At The National Oceanic and Atmospheric Administration (NOAA), containers provide an easy way to build and distribute complex HPC applications, allowing faster collaboration, portability, and experiment computer environment reproducibility amongst the scientific community. The challenge arises when applications rely on Message Passing Interface (MPI). This necessitates investigation into how to properly run these applications with their own unique requirements and produce performance on par with native runs. We investigate the MPI performance for benchmarks and containerized climate models for various containers covering selection of compiler and MPI library combinations from the Cray provided Programming Environments on the Cray XC supercomputer GAEA. Performance from the benchmarks and the climate models shows that for the most part containerized applications perform on par with the natively built applications when the system optimized Cray MPICH libraries are bound into the container, and the hybrid model containers have poor performance in comparison. We also describe several challenges and our solutions in running these containers, particularly challenges with heterogeneous jobs for the containerized model runs.

Index Terms—containers, climate, benchmark, cray

I. INTRODUCTION

With the advent of containers in the form of LXC [1], they have been continuously considered and evaluated [2] for use in High Performance Computing (HPC) environments. The creation of Docker brought containers into the mainstream

and further increased its demand in various applications. Containers found great use in cloud deployments as an easy way to package applications and create reproducible builds in a clean and easily distributable form. The advantages provided by the ease of packaging and reproducible builds cannot be understated, especially in HPC. Projects like Spack [3] were created to address this gap in easily building and packaging applications of libraries in HPC and have found great success, seeing wide use in many HPC centers. However, ease of distribution and reproducibility was still a standing problem that needs to be addressed, especially for providing standalone environments for individual users. The use of container technologies like Singularity aim to fill that gap in HPC, by offering the ability to create a single Singularity Image Format (SIF) file that can be passed around from developer to user in a center, or whose build recipe can be adapted for other centers with different systems.

The National Oceanic and Atmospheric Administration (NOAA) develops and runs many environmental computer models that run for various environmental systems at many different time scales. NOAA's Geophysical Fluid Dynamics Laboratory (GFDL) develops weather and climate models which require a tremendous amount of computing resources and must adhere to strict scientific reproducibility guidelines. NOAA's supercomputer GAEA [4], a HPE Cray XC40 system was purpose built to develop and run these models at scale to serve the needs of climate prediction. However, these models built natively rely on the specific software stack available on the system itself, so reproducing these models effectively at other centers can become a challenge, especially when reproducibility depends on old versions of compilers and dependencies such as netCDF. Containers solve this problem by being able to package the software environment needed

This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of Energy (DOE). The publisher acknowledges the US government license to provide public access under the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

for the model in a reproducible artifact that can then be distributed to researchers and users who do not have access to the GAEA system, and potentially be distributed to users on other systems. Users can build and run their own containers without the struggle of trying to exactly match the bespoke software environment that the models rely on.

However, using these containerized environments to run these models can be a challenge as these models as well as many other HPC applications rely on MPI for parallel computation. Evaluating the compatibility of MPI between the container and the host, as well as performance of containers with MPI under different models becomes a necessity. It is also important to look at how well the stated ABI compatibility between different MPI vendors work in practice for our containerized applications when we bind the host MPI libraries into the container where the application was built with a different MPI library. The fully coupled climate models rely on a resource manager's ability to run heterogeneous jobs or heterogeneous job steps with a shared `MPI_COMM_WORLD` and have as such been tailored to work with GAEA's Slurm resource manager's ability to run heterogeneous jobs. The use of heterogeneous jobs is not something that has been evaluated in previous work on studying containers for HPC. We will be looking at the issues that arise when trying to run such containerized models in heterogeneous jobs alongside studying the MPI compatibility and performance of the models with containers for heterogeneous jobs. Understanding and evaluating these issues would provide great benefit to the GAEA user base who can take prebuilt containerized models with the confidence of knowing the level of performance they will be getting out of these container images.

In Section II we cover some background on using MPI with containers and the different models available for that. In Section III we look at past work that examine this area of container performance under different conditions on HPC systems and identify the gap in evaluating containerized climate models at scale and containerized heterogeneous jobs. In Section IV we describe the climate models we use in this study. In Section V we describe the containers and host environments, the experiments we performed with benchmarks and climate models, and describe the results. In Section VI we discuss the challenges in getting to native like performance for the containers and discoveries made for the getting heterogeneous jobs working for the containerized climate models. In Section VII we discuss potential future paths we could explore as we move forward in standardizing the use and distribution containers for GAEA and other systems.

II. BACKGROUND

Generally speaking, as described by LXC [1], containers are a lightweight virtualization method with the core idea of creating a clean operating system environment without the need for an independent kernel. Today, projects looking to use containers have access to several different container platforms. Approaching HPC with the idea of using containers comes with its limitations though, in terms of aligning the container

platform with the HPC facility. The details of which we won't cover here, but is accessible in other texts, like Younge et al [5].

Our intention is to evaluate and exercise the portability and performance experience for selected containerized NOAA climate model applications. This is a tricky problem with and without containers. The idea of having an upstream repository of containers that can just run on a wide array of systems is very appealing though. The issue is the diversity of HPC systems, across different centers and under different vendors. A core reason is the dependency on certain libraries, provided by the targeted HPC system, like those underlying MPI. What we explore focuses on that specific struggle. We analyze how to build and run containerized MPI applications while looking to achieve performance and portability.

Container builds follow one of two models to fulfill the MPI library dependency; either applying the hybrid or bind model for including MPI [6]. The hybrid model refers to the method of building and running the MPI-enabled container in a way that the MPI inside the container will work with the MPI outside the container. Doing this successfully requires a couple of things. First, it is required that the MPI inside the container is compatible with the MPI on the targeted system. Second, if performance is critical, it is up to the container developer to optimally build the MPI implementation for the hardware being targeted. The bind model revolves around the idea of integrating the system provided MPI directly into the container. You can do this by building the container directly on the system or building the container remotely with a temporary, properly compatible, MPI. The temporary MPI you built with remotely must be ABI compatible with the provided system MPI. The bind model gives you a clearer path to performance, since you can use the optimized system MPI, but it makes portability more difficult. In our case, we are testing Cray-MPICH and Intel-MPI from the host, and GNU-MPICH and Intel-MPI in the containers. This means we are able to take advantage of some work from the MPICH ABI compatibility initiative [7]. We can predetermine what is ABI compatible, build the container remotely, and enable the ability to grab the system MPI at runtime. In the rest of this paper we look at the results of doing this in practice, with real applications from NOAA, on a production Cray system.

III. PRIOR WORK

Many prior work have evaluated the overhead produced by containers on different platforms, including on HPC systems. Morabito [8] identified that they have negligible overhead even in extremely resource contained systems like Raspberry Pis, making containers suitable for a wide variety of large and small systems. Full virtualization is unsuitable for HPC because of its incredible overhead, but containers almost entirely eliminate that overhead while still providing a consistent repeatable environment that can be passed around to others. Xavier et al. [9] and Beserra et al. [2] compare the performance of earlier container runtimes like LXC [1] and hypervisor based virtualization on HPC systems to identify

that containers have negligible overhead and much better performance for IO than virtualization. Similarly, Torrez et al. [10] have performed a thorough study at HPC scale with 512 nodes to see how more modern HPC focused container solutions like Charliecloud, Shifter, and Singularity produce almost zero performance overhead on the system and had very little performance difference between them, giving freedom to HPC centers to provide whichever container framework is most suitable for its users. This provides a consensus to the findings of several others like Rudy et al. [11] who studied the performance overhead and scaling of a containerized multiphysics biological simulation of a pulsatile artery with Docker, Shifter, and Singularity and demonstrated Shifter and Singularity’s comparable superiority in HPC environments; with Younge et al. [5] who conduct a performance study of Singularity on a Cray XC system at Sandia and find little performance difference from native for MPI benchmarks when using the native Cray MPI libraries and compared performance with a deployment on Amazon EC2; and Le et al. [12] who did MPI performance comparisons using NEURON, OSU, and Intel MPI benchmarks between native and Singularity and found little difference in performance at the San Diego Supercomputing Center. Younge et al. [5] also test the dynamic linking of vendor MPI libraries into containers and its effect on performance, which we also investigate as part of our work. Ruiz et al. [13] discovered potential bottlenecks for containers used with network bound applications due to the default network settings, which could have implications for MPI applications in HPC. Abraham et al. [14] studied the impact of container usage on the Lustre storage system, and the impact of the storage system on the container application performance, and found that Singularity due to its single file structure of its image has the least impact on the MDSs and the OSSs, which further motivates our use of Singularity on GAEA which relies on Lustre.

IV. CLIMATE MODELS

Two of GFDL’s models have been selected for use in this study. The first is the aquaplanet model based on the physics of the GFDL Atmosphere Model 4 [15] [16]. The aquaplanet is an idealized full physics atmosphere model that runs without any land or ocean variation. Instead, the surface is an “ice model” that is smooth and has a constant surface forcing. This model is simpler than the full AM4, runs faster, and ensures that we are only testing the atmosphere model so that other model components are not contributing to any load imbalance. This model can be used for studying large scale atmospheric phenomenon such as the Walker Circulation [17] and global tropical cyclone frequency [18].

The second model is the fully coupled GFDL Earth System Model 4.1 (ESM4) [19]. This model runs with an updated AM4 atmosphere model with full chemistry and the GFDL LM4P model coupled with the Modular Ocean Model (MOM6) and the Sea Ice Simulator version 2 (SIS2) ice model. The AM4 and LM4P run on one set of MPI ranks while the MOM6 and SIS2 runs on a different set of MPI ranks.

Two separate communicators are set up for each of the sets of ranks referred to in the model as a “pelist”. Communication happens in the coupler portion of the model using the MPI_COMM_WORLD communicator between the different model components on the separate pelists. The method that the ESM4 runs on Gaea is being referred to as a heterogeneous run in this paper.

Both models are run for a month of simulation time from January 1-31. The month run test is used by GFDL to calculate the simulated years per day (SYPD): the number of simulated years that the model can run in an actual day on the computer. The times reported are the main loop times and ignores the initialization and finalization. The main loop is where the actual simulation occurs in the model, and the one month benchmark has been successfully used in the past to calculate the SYPD.

V. EVALUATION AND RESULTS

We ran our experiments, both small and large scale, on Gaea: a HPE Cray XC 40 system located in the Oak Ridge National Laboratory and operated for NOAA by the Department of Energy. For production, Gaea’s C4 cluster consists of 2656 nodes, with each node running 2 18-core Intel Xeon Broadwell CPUs and 64 GB of memory, with a peak performance of 3.52 petaflops, and makes use of Cray’s Aries/Dragonfly interconnect. For our small scale runs we used T4, a 20 node test and development system, with each node running 2 16-core Intel Xeon Haswell CPUs and also with 64 GB of memory. Both clusters have access to the center-wide DDN Lustre scratch file system with 38 petabytes of space.

For our evaluations of the MPI based programs, we executed runs natively built with the Intel and the GCC compilers, as well as using three container images with different software environments. Separate runs were performed for hybrid and bind MPI models for each container to get the full scope of performance. Table I describes the relevant software packages on the host and in the container images, for the benchmark and climate model applications we used. We note that the climate model experiments primarily used the Intel compiler, and the ESM4 model and aquaplanet model with GNU compilers did not work with the hybrid model.

A. OSU Benchmark runs

As an initial performance evaluation to see how well the containers performed under heavy MPI communication workloads, we ran the `osu_allgather` benchmark from the OSU Micro Benchmark suite [20]. We ran separately compiled benchmarks for the Host Environment as well as separately compiled in each container, so that we are not mixing binaries built for different environments, host or container. For each container, we ran the benchmark using both the hybrid and the bind MPI model. For the GNU MPICH container, we ran the bind model tests by binding the Cray MPICH (ABI compatible) libraries built for GCC, as well as binding the non-Cray Intel MPI libraries just to see how far the ABI compatibility could stretch and see how their performance

TABLE I: Container images and the package versions used

Environment	Packages
GNU MPICH container (hpcme_gnu_rhel8.sif)	RHEL 8.4 GCC 8.4.1 MPICH 3.3.2 Netcdf 4.8.0
Intel MPI container (ubuntu-intel_2022.1.1.sif)	Ubuntu 18.04 Intel 2021.5 Intel MPI 2019.10 GCC 7.5.0 Netcdf 4.8.0
Intel 2021 OneAPI container (intel2021.2_netcdfc4.7.4_ubuntu.sif)	Ubuntu 18.04 Intel 2021.2 Intel MPI 2021.2 GCC 7.5.0 Netcdf 4.7.4
Host Environment	Intel 19.0.5 (PrgEnv-intel) GCC 8.3.0 (PrgEnv-gnu) Cray MPICH 7.7.11 Intel MPI 2019.5

differed. For the Intel containers, we only bind mounted the Cray MPICH libraries built for Intel for the bind model runs. We ran each of these on the small scale on the T4 cluster, with each run on 16 nodes running 32 ranks per node. The results of the small scale runs were taken as the median of 5 repetitions. The large scale runs on C4 used 450 nodes with 32 ranks per node but due to time limits were only able to do 1 repetition.

With the small scale runs and also in the large scale runs, we can already see the clear difference between hybrid and bind models in time to complete the `osu_allgather` benchmark. There is an order of magnitude difference between the Host Cray MPICH bind containers and the other containers (the hybrid and the GNU MPICH-Host Intel MPI bind container) on allgather performance for message sizes up to 8 bytes. The Host Cray MPICH bind containers are close to or on par with the native runs in performance for the small message sizes, with the GNU MPICH-Host Cray MPICH bind container being the slowest among them. Figure 1a and Figure 1b shows this difference for the small and large scale runs, without the GNU MPICH hybrid container. The GNU MPICH hybrid container in the large scale run performs the worst and jumps several orders of magnitude worse starting at 8 bytes. Figure 2 shows the comparison for large scale runs up to message size of 8 bytes on log scale in order to display the GNU MPICH hybrid performance. Further analyses of these benchmarks will not consider the GNU MPICH hybrid container.

Another interesting thing to note about the large scale runs is that the GNU MPICH container-Host Intel MPI bind run, the Intel 2021 OneAPI container hybrid run, and the Intel MPI container container hybrid run display these spikes in their time to complete at the 16 and 64 byte message sizes during the large scale run. Figure 3 shows the graph where the spikes are visible in comparison with the other runs. This could be a quirk of the GAEA networking or the PMI communication or something to do with how the MPI parameters on GAEA are tuned. Further investigation is required to make sure actual containerized application that require those message sizes

TABLE II: GNU MPICH-HOST Intel MPI bind vs Intel 2021 OneAPI-Host Cray MPICH bind upto 512 byte message size

Message Size	GNU MPICH-Host Intel MPI bind	Intel 2021 OneAPI-Host Cray MPICH bind
1.00	2106.26	172.13
2.00	3705.07	266.33
4.00	5357.16	363.31
8.00	5982.68	612.30
16.00	311939.36	1022.08
32.00	5132.17	1663.66
64.00	38318.42	2983.23
128.00	7422.58	5457.56
256.00	11731.34	10497.98
512.00	18970.68	19832.21

aren't affected and to determine the cause of the spikes.

For the large scale runs, the GNU MPICH container hybrid model run and the GNU MPICH container-Host Intel MPI bind model run could only complete up to 512 bytes message size before the job timed out. There seems to be some setup time overhead for the GNU MPICH-Intel MPI bind run that is not being measured in the benchmark itself because the actual measured performance gets closer to the other bind model runs at 256 and 512 byte message size. Table II shows the the performance from 1 to 512 byte message sizes comparing the GNU MPICH-Host Intel MPI bind container and the Intel 2021 OneAPI-Host Cray MPICH bind container.

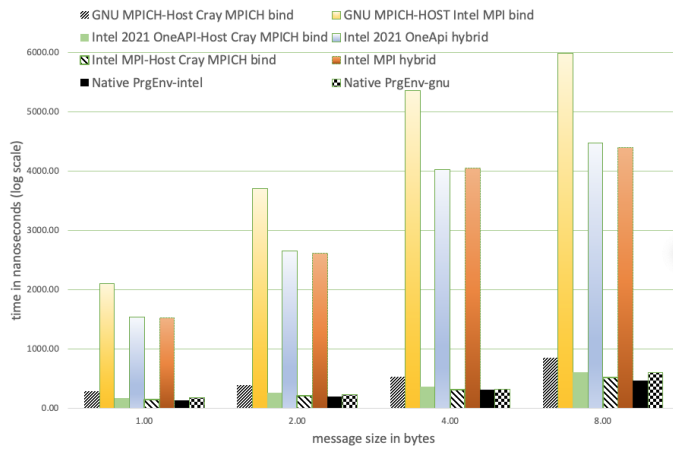
In the large scale runs, there are a couple of instances at 8192 and 16584 bytes message size where the Intel container hybrid runs seem to perform better than the native or the Intel container bind runs. Figure 4 shows this. Further investigation could net useful information on why this might be happening. But from an overall perspective, it is clear that the Intel containers with the Cray MPICH libraries bind mounted are the closest to native performance, and indicates to us that we should target that set up for containerizing applications to run on GAEA.

B. Aquaplanet model runs

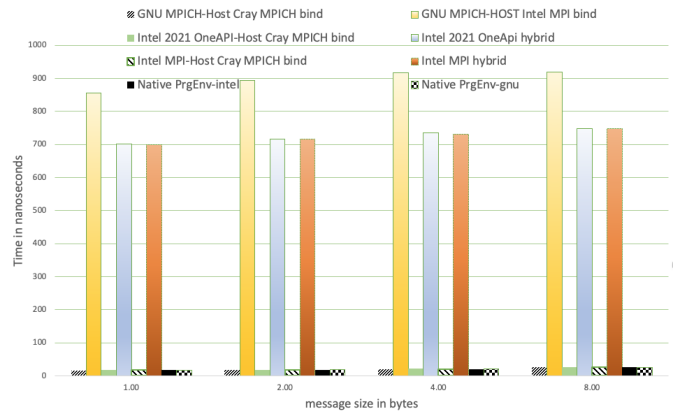
The aquaplanet model was scaled using an increasing number of MPI ranks and run on bare metal, using a Singularity Hybrid-MPI method, and the Singularity MPI-Bind method. The container and executable for both hybrid and bind are identical. The aquaplanet models were compiled using the intel oneAPI 2021.2 compilers and gfortran/gcc 8.4 for the intel and GNU runs, respectively. The container builds use a newer version of netCDF than bare metal.

The timings for the month runs of the aquaplanet with intel are found in Table III. The hybrid model runs much slower than the bare metal runs. As the number of cores increases, the percent difference increases from 12% for 216 cores to 102% for 1728 cores. For smaller and short runs, such as for development purposes, the hybrid model would be a viable choice, but for production runs at scale, the hybrid method is not a reasonable methodology.

The MPI-Bind timings are similar to the bare metal runs. The difference between the bare metal and MPI-Bind timings



(a) Large scale - 450 nodes, 32 tasks per node



(b) Small scale - 16 nodes, 32 tasks per node

Fig. 1: Performance of osu_allgather for message sizes up to 8 bytes

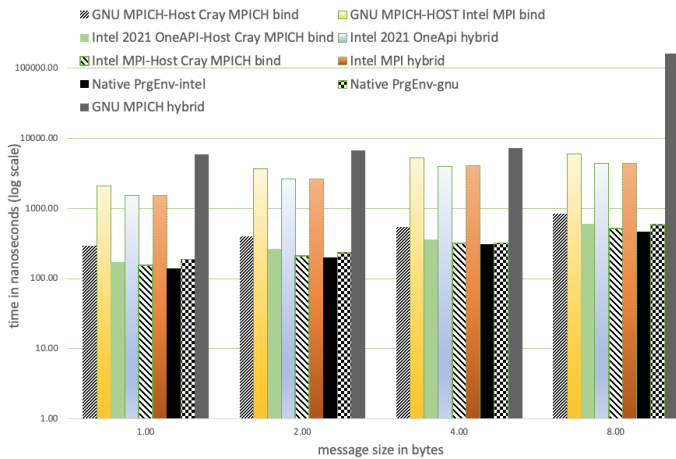


Fig. 2: Performance of large scale run of osu_allgather (450 nodes, 32 tasks per node) for small message sizes in log scale

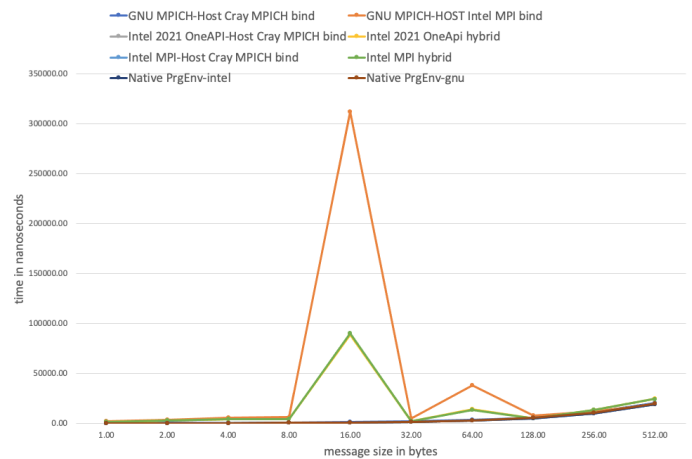


Fig. 3: Performance of large scale runs up to a message size of 512 bytes. Spikes observed at 16 and 64 byte message size

TABLE III: Number of cores (column 1) with main loop timings (in seconds) from one month runs of the aquaplanet model using bare metal, and the Singularity MPI-Hybrid and Bind configurations in the intel2021.2_netcd4.7.4_ubuntu.sif container.

total cores	bare metal (s)	hybrid (s)	bind (s)
216	1533.07	1745.50	1587.56
384	1194.56	1467.48	1217.80
432	847.57	1131.60	850.40
576	632.96	946.75	635.71
864	467.62	722.46	464.47
1152	378.64	775.05	392.15
1728	283.91	572.44	302.46

is within the difference in timings experienced on a run-to-run basis on the machine. Compared to multiple bare metal runs (not shown), the MPI-Bind timings are faster in some cases. While there is a steeper learning curve and a great deal of system work to be done in order to figure out what needs

TABLE IV: Number of cores (column 1) with main loop timings (in seconds) from a one month run of the aquaplanet model using bare metal, and the Singularity MPI-Bind configurations in the hpcme_gnu_rhel8.sif container.

cores	bare-metal (s)	bind (s)
216	2026.35	1951.29
384	1597.43	1532.81
432	1108.32	1103.46
864	877.39	601.98

to be bound into the container to get the MPI-Bind model to work, the MPI-Bind method is a viable use case for production runs.

Similar results are observed in Table IV when the aquaplanet was run in MPI-Bind mode using the GNU container. The timings are all within the natural variation of run-to-run differences, except for the run with 864 cores where the container was about 30% faster. The bare metal run appears to

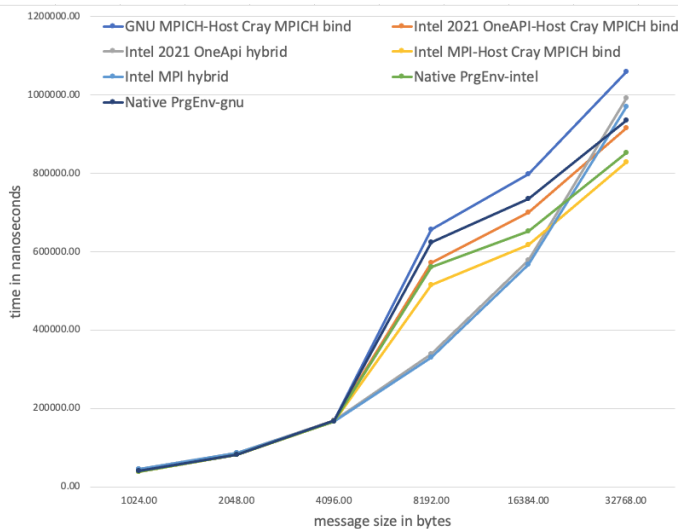


Fig. 4: Performance of large scale runs for large message sizes. Hybrid runs of the Intel containers seem to perform better in a couple of instances.

be losing its scaling, but the container seems to maintain the scalability. This result would become important for a higher-resolution model that might be run on many thousands of cores, such as a weather model. It may be possible to get better scaling performance using a container based on this result.

C. ESM4 runs

The ESM4 model programs were compiled in the same manner as the aquaplanet model. Then the standard one-month timing test suite was run on both the Host and Container environments. Table V shows the test experiment configurations; the coupled ESM4 model runs the atmosphere/land and ocean/ice on separate sets of MPI ranks, and OpenMP threads are used for the atmosphere and not the ocean model component. These configurations are chosen so that the atmosphere and ocean components run in a similar amount of time.

Challenges were overcome in order to run the coupled ESM4 model within the container environment, including binding needed library directories, setting specific environment variables, and including the submitting Gaea job environment in the batch environments. Additionally, only the MPI-Bind method was used, as the container’s Intel MPI is currently incompatible with heterogeneous jobs. These are described more in the Discussion section.

The results between the Native Intel and the MPI-Bind timing experiments are shown in Table V. While the MPI-Bind container runtimes were overall 1-2% slower on average, the differences were within the natural run-to-run variation. Overall, the results are consistent with the OSU benchmark and aquaplanet experiments in showing similar performance between Native and MPI-Bind container applications.

While there is more work to understand and optimize which system resources and directories need to be bound in the container, it is encouraging that when using the MPI-Bind

TABLE V: ESM4 experiment configurations (nodes, atmosphere MPI ranks and OpenMP threads, and ocean MPI ranks) with main loop timings (in seconds) from a one month run of using Native and MPI-Bind configurations in the intel2021.2_netcdcf4.7.4_ubuntu.sif container. (Timings are an average of 2-3 runs. Ocean model does not use OpenMP.)

nodes	atmos ranks x threads	ocean ranks	Native (s)	Container (s)
30	216 x 2	648	2898.5	2899.6
36	216 x 2	951	2904.5	2971.4
51	432 x 4	951	1675.0	1732.4
54	432 x 4	1057	1720.7	1734.5
60	432 x 4	1296	1685.3	1732.9
61	432 x 4	1300	1656.1	1714.9
69	576 x 4	1300	1366.0	1417.9
85	864 x 4	1300	1057.8	1047.0
89	864 x 4	1441	1018.8	1016.0
93	864 x 4	1608	975.8	1010.6
94	1152 x 4	1057	1072.5	1065.0
105	864 x 4	2044	971.5	1010.7

method, the performance impact of containerization on real-world MPI applications such as the ESM4 model is low.

VI. DISCUSSION

A useful feature of containerization is the ability to simply package up the application and the required dependencies, including any MPI libraries, in a single container image to pass on to other users. However, from our evaluation there is a clear benefit in being able to bind the optimized host Cray MPICH libraries into the container, provided the host MPI libraries are ABI compatible with the container’s MPI libraries. This is useful to know for any given supercomputer so that the appropriate paths can be mounted into the container. However, the process of binding these paths is manual where the appropriate MPI module, that is ABI compatible with the MPI library in the container, has to be loaded if it even exists on the system, and then add the appropriate paths pointing to the `libmpi.so` location as mounts in the container with the `--bind` flag. This can be a process of trial and error when it comes to Cray MPICH installations because in addition to the `libmpi.so`, one also needs to make sure that all the supporting libraries that `libmpi.so` depends on are also mounted in the container. This includes Cray specific libraries like `alps`, `ugni`, `udreg`, `xpemem` which provide host specific optimizations for MPI communications. Part of the reason for the slowness of the hybrid model runs may come from the container’s internal MPI libraries not having access to these optimizations. There is no standard documentation for mounting these for container runs so figuring it out took some time and effort. Being able to automatically discover these libraries and mount them would be a useful feature in our system and also useful if the containers we build are passed to other users elsewhere who just want to be able to run it with maximum performance without needing to go through the excess effort in finding all the host libraries they need. `e4s-cl` [21] was built for this, where it is able to automatically discover the MPI libraries and creates a profile of the libraries to mount into the container. A container run using `e4s-cl` will

automatically bind the discovered libraries that are saved in the specified profile into the container. However, this is currently still a work in progress and the automatic discovery is limited since it requires the use of mpicc as part of the discovery process, which GAEA does not have.

The need to bind mount the host's MPI libraries into the container in order to get significantly better performance does break some portability. It is not so simple as just taking a container image and running. It is possible with future improvements to Cray's PMI could bring performance of the hybrid model runs of containers closer to that of the bind model and native runs, which would make usage of containers on HPC so much more seamless and save users the hassle of needing to manage ABI compatible libraries for binding.

The ESM4 model's reliance on heterogeneous jobs exposed issues in both host and container runtime environments. Slurm allows both heterogeneous job allocations and job steps. Running jobs with the host MPI library resulted in hangs on `MPI_Init` with heterogeneous allocations, but ran successfully with heterogeneous job steps. Intel's MPI library was incapable of receiving the heterogeneous task distribution, resulting in incorrect creation and access attempts of shared memory resources. The MPICH library in GNU toolchain container ran successfully only if the job allocation was made with the `hetjob` flag, and hung on `MPI_Init` otherwise. Updating Slurm on a similar test system has resolved the host MPI library hang on `MPI_Init` while the other libraries maintain the same behavior. This has made the linking of the host MPI library into the container's executable important for both performance and basic functionality.

When running performance tests for containers with Cray MPI libraries, it was initially found that the container with the GNU toolchain did not benefit from the addition of the library path to `LD_LIBRARY_PATH`. This was due to the container's `mpicc` compiler wrapper setting `RPATH` for the container's MPI library. Multiple workarounds were successfully tested including compilation with `-Wl,--enable-new-dtags` to instead set `RUNPATH`, setting `LD_PRELOAD` for the Cray MPI library, and launching the executable in the container by invoking the dynamic linker directly with the `--inhibit-rpath` flag. Each method was found to have its drawbacks. For the first, the user of the container may not have the means to rebuild the application. For the second, loading the library through the environment variable applies to a wider scope than just the executable we wish to run. In the last case, explicitly excluding search for MPI libraries in `RPATH` when a matching host library cannot be found can result in several different outcomes. Since version 2.30, the GNU dynamic linker offers a `--preload` flag to preload a library only for the specified executable. When available within the container this option offers the ability to reliably load host libraries in the proper scope.

The Slurm scheduler (used on GAEA and other HPC sites) allows users to specify which environment variables are passed from the submitting environment (login/head node) to the batch environment (compute nodes). While the default is

to pass all environment variables through `--export=ALL`, the most common use on Gaea is to pass no environment variables through `--export=NONE`, which then populates the batch job environment using the system default shell initialization. This helps ensure a clean batch environment for reproducible experiments. However, we discovered that using `--export=NONE` caused heterogeneous jobs (simple tests as well as ESM4) to hang immediately after launching. More work is needed to discover which environment variables in the login user environment are needed for heterogeneous jobs.

Containerizing the climate models can be time consuming, but should not be detrimental to the overall development of climate model code. The most time consuming part of dealing with climate model containers is compiling packages and dependencies that are required, such as HDF5 and netCDF. There is currently an effort to cut down this overhead to make containers even easier to build. Also, providing containers that can be used with the Intel compiler poses unique licensing challenges because there is not clarity as to whether we can share containers that have the compilers necessary to compile the model. Once these hurdles are overcome, there were only slight modifications to the GFDL workflow that were needed to compile and run the model. The portability and shareability of the container make it an attractive choice for the future of climate modeling.

VII. FUTURE WORK

There are some further things we could evaluate as we move forward in getting containers ready for regular use on GAEA. Testing how well the Lustre filesystem handles reads from the container images for very large scale deployments is of importance in the exascale era especially if the SIF files are large and aren't properly striped. Automatically discovering and bind mounting the right libraries from the host into the container is another useful area to look at. `e4s-cl` [21] has done some work in this regard but currently relies on the presence of `mpicc` and `mpirun`, which does not work on a Cray system if it doesn't have those executables and may not correctly pull in the necessary Cray libraries. Exploring faster builds for model software dependencies packages is planned to be done through a partnership with E4S utilizing spack caches and large base containers. Using the containers to run on other HPC platforms and the cloud is a planned goal for NOAA. The next-generation models will likely be shared with containers instead of just a repository on GitHub.

VIII. CONCLUSION

In this paper, we explored the containerization of MPI based applications and climate models on the GAEA supercomputer then evaluated their performance under different conditions. We identified that binding the host's MPI libraries into the container nets the best performance when done correctly, as well as identifying some odd behavior during benchmark runs for some hybrid model containers. We explored the challenges and figured out solutions for containerizing and running applications that specifically rely on the ability to run

heterogeneous jobs. We articulate steps needed to effectively run those applications at native performance without needing to modify the applications themselves to suit the more conventional non-heterogeneous job blueprint. We discuss the future benefits of containerization for distributing climate models and engendering collaboration in the climate science community.

ACKNOWLEDGMENT

This research used resources of the National Climate-Computing Research Center which is located within the National Center for Computational Sciences at the Oak Ridge National Laboratory and supported under a Strategic Partnership Project between the DOE and NOAA. This research used resources of the Oak Ridge Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC05-00OR22725.

REFERENCES

- [1] "LXC." <https://linuxcontainers.org/lxc/>. Accessed: November 25 2019.
- [2] D. Beserra, E. D. Moreno, P. T. Endo, J. Barreto, D. Sadok, and S. Fernandes, "Performance analysis of lxc for hpc environments," in *2015 Ninth International Conference on Complex, Intelligent, and Software Intensive Systems*, pp. 358–363, IEEE, 2015.
- [3] T. Gamblin, M. LeGendre, M. R. Collette, G. L. Lee, A. Moody, B. R. De Supinski, and S. Futral, "The spack package manager: bringing order to hpc software chaos," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–12, 2015.
- [4] "GAEA Supercomputer." <https://www.noaa.gov/organization/information-technology/gaea>. Accessed: April 14 2022.
- [5] A. J. Younge, K. Pedretti, R. E. Grant, and R. Brightwell, "A tale of two systems: Using containers to deploy hpc applications on supercomputers and clouds," in *Proceedings of the 2017 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, Dec. 2017.
- [6] "Singularity and MPI." <https://sylabs.io/guides/3.5/user-guide/mpi.html>. Accessed: April 14 2022.
- [7] "MPICH ABI Compatibility Initiative." https://wiki.mpich.org/mpich/index.php/ABI_Compatibility_Initiative. Accessed: April 14 2022.
- [8] R. Morabito, "A performance evaluation of container technologies on internet of things devices," in *2016 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, pp. 999–1000, IEEE, 2016.
- [9] M. G. Xavier, M. V. Neves, F. D. Rossi, T. C. Ferreto, T. Lange, and C. A. De Rose, "Performance evaluation of container-based virtualization for high performance computing environments," in *Proceedings of the 21st Euromicro International Conference on Parallel, Distributed, and Network-Based Processing.*, IEEE, 2013.
- [10] A. Torrez, T. Randles, and R. Priedhorsky, "Hpc container runtimes have minimal or no performance impact," in *2019 IEEE/ACM International Workshop on Containers and New Orchestration Paradigms for Isolated Environments in HPC (CANOPIE-HPC)*, pp. 37–42, 2019.
- [11] O. Rudyy, M. Garcia-Gasulla, F. Mantovani, A. Santiago, R. Sirvent, and M. Vázquez, "Containers in HPC: A Scalability and Portability Study in Production Biological Simulations," in *Proceedings of 2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, May 2019.
- [12] E. Le and D. Paz, "Performance Analysis of Applications using Singularity Container on SDSC Comet," in *Proceedings of the Practice and Experience in Advanced Research Computing 2017 on Sustainability, Success and Impact (PEARC).*, (New Orleans, LA, USA), 2017.
- [13] C. Ruiz, E. Jeanvoine, and L. Nussbaum, "Performance Evaluation of Containers for HPC," in *Proceedings of the 2015 European Conference on Parallel Processing*, 2015.
- [14] S. Abraham, A. K. Paul, R. I. S. Khan, and A. R. Butt, "On the use of containers in high performance computing environments," in *2020 IEEE 13th International Conference on Cloud Computing (CLOUD)*, pp. 284–293, IEEE, 2020.
- [15] M. Zhao, J.-C. Golaz, I. Held, H. Guo, V. Balaji, R. Benson, J.-H. Chen, X. Chen, L. Donner, J. Dunne, *et al.*, "The gfdl global atmosphere and land model am4. 0/lm4. 0: 1. simulation characteristics with prescribed ssts," *Journal of Advances in Modeling Earth Systems*, vol. 10, no. 3, pp. 691–734, 2018.
- [16] M. Zhao, J.-C. Golaz, I. Held, H. Guo, V. Balaji, R. Benson, J.-H. Chen, X. Chen, L. Donner, J. Dunne, *et al.*, "The gfdl global atmosphere and land model am4. 0/lm4. 0: 2. model description, sensitivity studies, and tuning strategies," *Journal of Advances in Modeling Earth Systems*, vol. 10, no. 3, pp. 735–769, 2018.
- [17] L. G. Silvers and T. Robinson, "Clouds and radiation in a mock-walker circulation," *Journal of Advances in Modeling Earth Systems*, vol. 13, feb 2021.
- [18] A. C. Burnett, A. Sheshadri, L. G. Silvers, and T. Robinson, "Tropical cyclone frequency under varying SSTs in aquaplanet simulations," *Geophysical Research Letters*, vol. 48, mar 2021.
- [19] J. Dunne, L. Horowitz, A. Adcroft, P. Ginoux, I. Held, J. John, J. Krasting, S. Malyshev, V. Naik, F. Paulot, *et al.*, "The gfdl earth system model version 4.1 (gfdl-esm 4.1): Overall coupled model description and simulation characteristics," *Journal of Advances in Modeling Earth Systems*, vol. 12, no. 11, p. e2019MS002015, 2020.
- [20] "Osu micro benchmarks." <https://mvapich.cse.ohio-state.edu/benchmarks/>. Accessed: March 28 2022.
- [21] "E4S-CL." <https://github.com/E4S-Project/e4s-cl>. Accessed: April 14 2022.