

# Unified Model Global Weather Forecast Performance on HPE Cray EX

Peter Johnsen

Hewlett Packard Enterprise  
Performance Engineering  
Bloomington, MN USA  
peter.johnsen@hpe.com

Steven Warren

Hewlett Packard Enterprise  
Performance Engineering  
Bloomington, MN, USA  
steven.warren@hpe.com

**Abstract**— The next generation HPE Cray EX (formerly Cray Shasta) supercomputer offers excellent performance for a wide range of applications including numerical weather prediction. In a compact architecture that includes AMD EPYC Rome processors along with the latest HPE Slingshot high-speed interconnect, the HPE Cray EX system is showing superb weather simulation performance. In this paper we look at the performance of the Unified Model (UM) from the UK Met Office. The UM is currently producing global and regional forecasts at a number of operational weather centers around the world. The UM global weather forecast ensembles at 10 km resolution are achieving net simulation speeds of up to 45 forecast days per wall clock hour on 700 nodes of the HPE Cray EX. This includes full model forecast output and shows very little run time variability across ensemble copies.

HPE Slingshot interconnect congestion management features and the impacts on UM performance while using the GPCNeT network load program are also investigated.

**Keywords**— component; UM, I/O, GPCNeT, network congestion, supercomputer, NWP

## I. INTRODUCTION

HPE Cray EX supercomputers are under deployment at a number of sites including operational weather forecast centers. We present the performance of the Unified Model (UM) global forecast model from the UK Met Office on the

HPE Cray EX systems currently under installation at Oak Ridge National Lab (ORNL) for the US Air Force 557<sup>th</sup> Weather (AFW) Wing.

A key performance aspect for any operational weather forecast center is the reliability of model run times, which ensures that strict production schedules can be met. This is required even under unpredictable and varied workloads. If one portion of a forecast cycle lags for any reason, including compute system load, final forecast products will not be available to consumers at expected intervals. For civilian and military forecasters alike, this can have serious consequences for critical weather situations.

The HPE Cray EX systems incorporate next generation technologies, including advanced high-speed networks (HSNs) which incorporate industry-leading adaptive routing and novel congestion management ensuring consistent run time performance. The HPE Slingshot [1] network is key to achieving fast and reliable weather model performance.

## II. HPE CRAY EX OVERVIEW

A pair of HPE Cray EX supercomputers [2] [3], along with a separate HPE Cray EX login and management server and DDN storage with Lustre file systems is being installed at Oak Ridge National Laboratory for US Air Force Weather forecasters. Each identical compute server (Systems 1 and 2) is comprised of identical components. Both systems have four

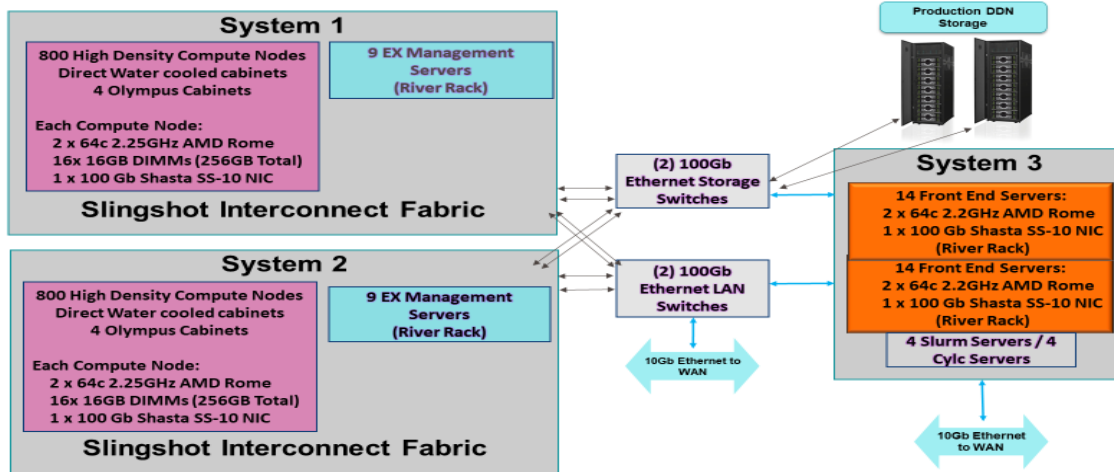


Figure 1 HPE Cray EX AFW System Design

high-density, liquid-cooled cabinets hosting AMD EPYC 7742 “Rome” processors (one switch group per cabinet). There are 200 dual-socket nodes per high-density cabinet consisting of 8, 12, or 16 nodes per Rosetta switch. Each cabinet is partially populated with compute nodes leaving room for future expansion with GPU accelerators for AI/ML applications. See Figure 1 and Table 1 for more details of the AFW HPE Cray EX systems.

Table 1: AFW HPE Cray EX Compute System Components

4 Olympus liquid cooled cabinets per system
800 nodes per system
2 AMD EPYC 7742 Rome 64 core processors (128 total cores) per node
256 GB DDR4-3200 memory per node
HPE Slingshot 10-100 Gbit/s interconnect, 32 Rosetta switches per system (5 groups)
Single NIC per node

### III. CRAY XC40 OVERVIEW

We also ran on a Cray XC40 platform to compare the performance differences of the UM between the previous generation of network and processor technologies with the current HPE Cray EX described in Section A. The Cray XC40 we used consisted of 6 cabinets (3 switch groups) of 484 dual-socket 18 core Intel Broadwell 2.1 Ghz processors. The Aries network is configured in a similar dragonfly topology to the HPE Cray EX machines. The Aries network also employs adaptive routing, responding to bottlenecks in network traffic, however, it lacks the advanced congestion management developed for the HPE Slingshot network and has been shown to be sensitive to certain congesting traffic patterns [4].

### IV. UM MODEL OVERVIEW

UK Met Office Unified Model is a comprehensive numerical weather prediction application and is the main global forecast model at not only the Met Office, but also other centers including the Australia Bureau of Meteorology, the Korea Meteorological Administration, and US Air Force Weather, among others. Other configurations of the UM are used for regional forecasting, data assimilation, and climate simulations. The model is available through agreement with the Met Office.

The UM’s ENDGame dynamical core uses a semi-implicit semi-Lagrangian formulation to solve the non-hydrostatic, fully compressible deep-atmosphere equations of motion. For a detailed description of the UM model, see [5].

Tables 2 and 3 list the specific UM model configuration used for our benchmark performance analysis along with the main HPE Cray EX system software components used. Note that the model resolution of 10 kilometers (grid spacing) compares with the 10-12km resolutions currently in production at most global forecast centers.

Table 2: UM global model configuration

UM version: 10.9, GA6.1
N1280 horizontal resolution ~10km
70 vertical levels
72 hour forecast
Output every forecast hour, 4.4 GBytes written each file
240 second time step, 1080 total time steps

Table 3: Programming Environment for code build

Cray CCE Fortran and C/C++ compiler version 10.0.1
Cray MPICH and libfabric version 8.0.11
Cray Scientific library (libsci) version 20.0.6.1.1
HPE Cray OS version 1.2.1

### V. UM THROUGHPUT PERFORMANCE

As part of the system performance requirements of the AFW HPE Cray EX installation, UM throughput tests were run on each of the two systems. The requirement was to nearly fill each system with five concurrent copies of the UM global model described in Table 2 achieving a committed level of performance while maintaining 5% or lower run time variability. The run time configuration for each copy is given in Table 4.

Table 4: UM run time configuration for each copy

135 nodes, 17,280 cores
Horizontal MPI decomposition 40x88 (4 OpenMP threads per rank)
8 IO server groups, 12 ranks per server group
Hyperthreads not used
Output every forecast hour, 4.4 GBytes written per file

Performance of the five concurrent copies shows a run time variation of only 1.2% despite using 675 nodes on the system. Run times in wall clock seconds as measured internally by the UM code are as follows:

```
out:um_shell: Info: End model run at time= 1179.433 seconds
out:um_shell: Info: End model run at time= 1185.882 seconds
out:um_shell: Info: End model run at time= 1185.605 seconds
out:um_shell: Info: End model run at time= 1192.940 seconds
out:um_shell: Info: End model run at time= 1178.277 seconds
```

Expressed in terms of simulation speed, each copy produces a 9-day forecast in one wall clock hour. Combined, that comes to a net 45 forecast days per wall clock hour.

Forecast output is a large part of any current weather forecast application and can drastically increase the time-to-resolution if not managed well. UM has built-in dedicated IO server tasks that can largely hide this IO overhead. When forecast output is scheduled, data is handed off asynchronously to the IO server ranks as fast as possible so that the computations can proceed. For our UM simulation, the time to initiate a data transfer is around 0.5 to 1 second

which adds up to only 75 seconds of total ‘stall’ time in the forecast computations. How we handle the IO handoff is an important part of our model configuration but effectively increases the amount of MPI communications. The increase of background MPI communication creates another potential source of interference in the messaging associated with the model computation, e.g., halo exchanges and global diagnostics. However, we see very little, if any, interference on the HPE Slingshot interconnect network, even when filling the system with multiple copies.

## VI. STRONG SCALING RESULTS

We gathered scaling performance on the HPE Cray EX system using the same configuration described in the previous section. Additional IO server ranks were required at the highest node count to manage IO server data collection traffic. As seen in Figure 2, a small degree of superlinear scaling is observed between 69 and 135 nodes due to better cache utilization at 135 nodes (and above), as well as expected performance differences in decomposition strategies across the node counts. At 270 nodes (34,560 cores), parallel efficiency is still above 80%.

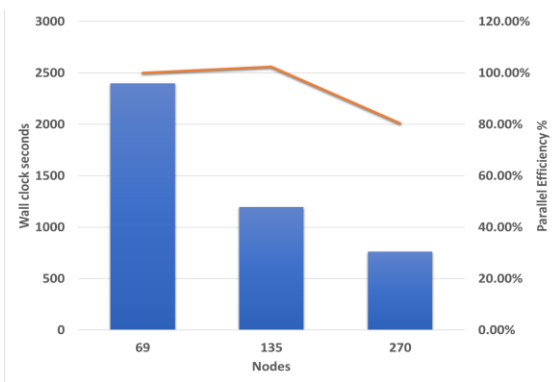


Figure 2: UM strong scaling results

## VII. UM COMMUNICATION ASPECTS AND NETWORK CONGESTION TESTING

The real challenge for any supercomputer is keeping pace with varied, and sometimes very aggressive, communication bound workloads. Applications that do very frequent small messaging, sometimes due to poorly constructed IO patterns, can cause congestion across high-speed networks resulting in unexpected variations in run times. As noted in our introduction, this is a critical problem for operational weather forecast centers and can cause a slowdown or worse in forecast product generation for its consumers.

In Figure 3 we show a high-level performance breakdown obtained through UM profiling. While the simulation writes significant output, IO is a small percentage of the overhead as seen from computation ranks. Communication overhead is not that large of a factor either, but as we will see later in this section, it is frequent enough and of the right type to make UM susceptible to network congestion from other sources.

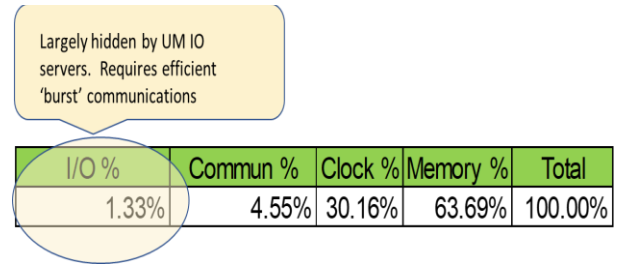


Figure 3: UM performance components

We have shown in Section C that there is very little network interference on the HPE Cray EX when running multiple copies of the UM. Another test for network congestion was performed on the AFW HPE Cray EX systems using a single copy of the UM on 135 nodes along with a much more aggressive network message generator, the GPCNeT benchmark [4]. GPCNeT is a freely available program designed to measure the impact of significant message rates and known communication patterns which cause congestion on high-speed networks. GPCNeT is composed of two sets of ranks, the first set being a ‘canary’ set that measures typical application MPI message performance, while the second set is a ‘congestor’ set that performs a large number of small but very frequent messages. For our tests, we minimized the canary set in favor of more congestors since we are interested in how this may impact the UM performance. We expect HPE Slingshot network congestion management and adaptive routine features [6] to limit interference with communications in the UM.

We ran the same tests on a Cray XC40 system, which uses the previous Cray Aries high-speed network and Intel Broadwell 18-core processors. Table 5 lists the test parameters for each system.

Table 5: Run configuration for combined UM and GPCnet tests

	Program	Nodes	MPI ranks	OpenMP threads	Total cores
<i>Cray XC40</i>	UM	360	4320	3	12960
	GPCnet	124	4464	1	4464
<i>HPE Cray EX</i>	UM	135	4320	4	17280
	GPCnet	640	81920	1	81920

The UM was first run standalone on each system to collect the best run time performance. The UM cases were mainly identical on both XC40 and HPE Cray EX systems, i.e., compiler, domain decomposition, IO servers, etc. with the following exceptions; 1) more nodes were needed on the XC40 since there are only 36 cores per node and 2) only 3 OpenMP threads per MPI rank were used on the XC40 to fit 18 core Broadwell processors. After the standalone UM runs had been completed, the UM was restarted and, after it had reached time step 50, GPCNeT was started on the remaining nodes.

The UM prints the wall clock time for each time step during the forecast. In Figure 4 the first 24 hours of the forecast (360 time steps) are plotted for each of the four tests.

GPCNeT ran from approximately UM time step 50 until time step 220. The steps with the highest wall times are those where the UM model is performing forecast output and include data offload to the UM IO server ranks.

Overall, we see that performance is better per time step on the HPE Cray EX than the XC40 due to using 4 OpenMP threads per rank rather than three and using the newer AMD Rome processors with a slightly faster base clock plus faster memory.

In Figure 4 we see clearly that the UM performance suffered by more than 30% on the XC40 when GPCNeT was running. Steps with and without forecast output were affected, output steps even more so. In contrast, no evidence of network congestion is apparent on the HPE Cray EX system even though GPCNeT was running on significantly more cores generating a higher number and total rate of messages.

The congestion management features of the HPE Slingshot interconnect allow the UM to achieve the same performance when it shares the system with other network-intensive workloads as it does when it runs standalone.

UM are run simultaneously. We observe no evidence of network interference when running the heavy messaging application GPCNeT alongside UM, highlighting the advantage the HPE Slingshot network's advanced congestion control plays in producing repeatable run times. The same UM model run on the Cray XC40 gets lower performance, larger run-to-run variability, as well as is shown to be susceptible to network load.

The HPE Cray EX solution with the HPE Slingshot high-speed network, with its advanced adaptive routing and congestion management technologies, will help operational weather forecast centers maintain rigorous production schedules despite often unpredictable workloads.

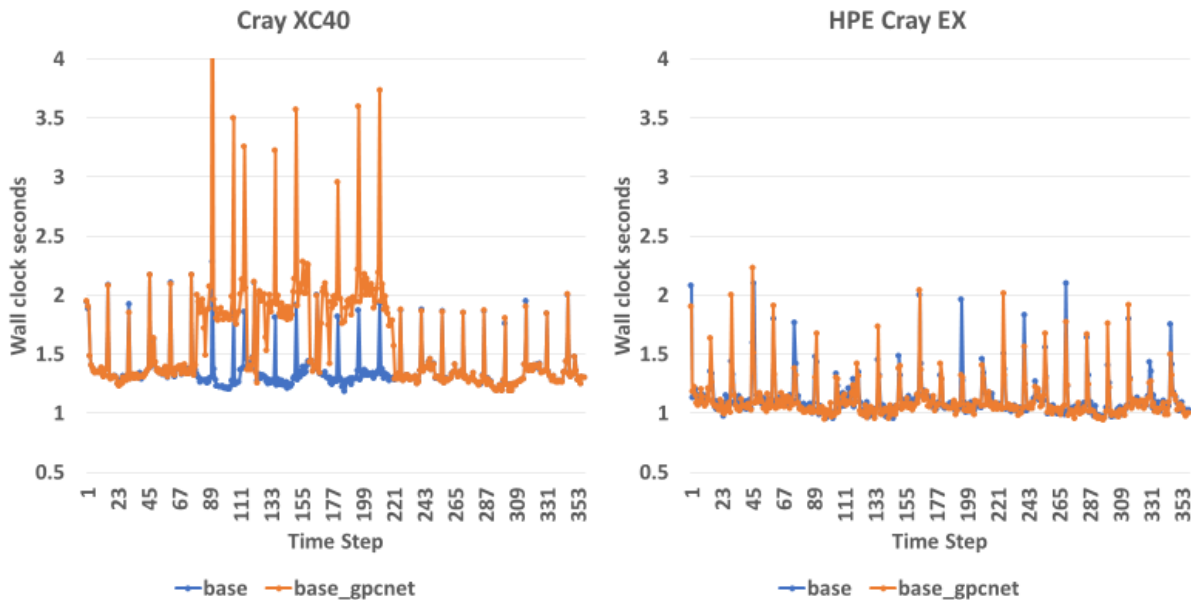


Figure 4: UM per time step duration for Cray XC40 (left) and HPE Cray EX (right). Blue 'base' is standalone run, orange 'base+gpcnet' is combined run

## VIII. CONCLUSION

In this paper we have shown the performance of the UK Met Office UM global forecast model on the new HPE Cray EX supercomputer system and compared to the previous Cray XC40 Aries-based system. UM model performance shows very little run time variation, even under heavy network loads, on the new HPE Cray EX platform. A small 1.2% wall clock variation is noted across copies when multiple instances of

## ACKNOWLEDGMENTS

We thank Oak Ridge National Laboratory and staff for the use of the HPE Cray EX systems during the installation period. We also thank the following HPE colleagues for the valuable comments and suggestions, Patricia Balle, Ilene Carpenter, and Chris Davidson.

## REFERENCES

- [1] HPE Cray Slingshot: Interconnect for the Exascale Era, <https://assets.ext.hpe.com/is/content/hpedam/a50002368enw>
- [2] HPE The Next Generation of Exascale Computing, <https://www.hpe.com/us/en/compute/hpc/supercomputing/Cray-exascale-supercomputer.html>
- [3] HPE Cray EX Supercomputer QuickSpecs, <https://h20195.www2.hpe.com/v2/GetDocument.aspx?docname=A00094635ENW>
- [4] Sudheer Chunduri, Taylor Groves, Peter Mendygral, Brian Austin, Jacob Balma, Krishna Kandalla, Kalyan Kumaran, Glenn Lockwood, Scott Parker, Steven Warren, and et al. Gpcnet: Designing a benchmark suite for inducing and measuring contention in hpc networks. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 19, New York, NY, USA, 2019. Association for Computing Machinery.
- [5] David Walters, Ian Boutle, Malcolm Brooks, and et al, The Met Office Unified Model Global Atmosphere 6.0/6.1 and JULES Global Land 6.0/6.1 configurations, Geosciences. Model Development., 10, 1487–1520, 2017
- [6] Daniele De Sensi, Salvatore Di Girolamo, Kim H. McMahon, Duncan Roweth, Torsten Hoefler, n In-Depth Analysis of the HPE Slingshot Interconnect, To be published in Proceedings of The International Conference for High Performance Computing Networking, Storage, and Analysis (SC '20) (2020)