

# **WRF performance on IaaS from the public cloud Fall 2020 (v1.0)**



**N**umerical Weather Prediction (NWP) is a fundamental asset in today's modern societies as weather forecasts are useful to citizens and businesses alike. Most present NWP run on on-premises HPC (High Performance Computing) hardware with less than 2% of current predictions being generated in cloud environments. However, shifting a larger fraction of these computations to the public cloud is desirable as cloud environments accelerate the adoption of new hardware and provide a framework to integrate computations, data assimilation and high-performance storage with nearly infinite capacity. Cloud Service Providers (CSPs) provide computational and storage options in the form of Infrastructure-as-a-Service (IaaS). The extent of these services can be augmented or lowered to fit fluctuating workloads or storage needs without being tied up to a monolithic configuration. From an economic perspective, the use of cloud resources eliminates the large upfront capital needed to purchase expensive equipment. This upfront capital and how to finance it are traditionally two of the major challenges for small and medium-sized businesses to access HPC capabilities. By circumventing this cycle, smaller organizations gain access to resources traditionally reserved to large businesses, national labs, and research universities. Two additional advantages from using IaaS are advanced security features and the reduction in maintenance costs as most of these tasks are performed by CSPs' technical teams.

Among the many NWP apps available nowadays, the Weather Research and Forecasting (WRF) Model is one of the most

popular ones with over 48,000 registered users from 160 countries. WRF is a mesoscopic model able to generate forecasts for a wide range of atmospheric conditions and resolutions. These forecasts serve for operational functions as well as for research purposes. The National Center for Atmospheric Research (NCAR) is the main developer of WRF, but several other national laboratories and academic institutions have contributed during its twenty plus years of development.

### **IaaS computational power**

CSPs offer computational power in the form of IaaS that can rival in performance to the supercomputers and clusters traditionally used to perform NWP. Even though early IaaS computational power used traditional server processors, most recent IaaS additions suitable for HPC workloads take advantage of custom processors specifically tailored to work in cloud environments. Some of these processors are variations of server processors but recent additions are designs specific for cloud environments. Examples of cloud exclusive processors are AWS Graviton 2 and the new Ampere® Altra™ processor. AWS introduced Graviton 2 at re: Invent 2019 and Ampere® Altra™ is currently in its latest phases of testing. Altra™ is expected to become available in the last quarter of 2020. The benchmarks with Altra™ presented here are therefore performed with preproduction hardware.

The present benchmarks use IaaS from 3 major CSPs: Amazon Web Services or AWS (<https://aws.amazon.com>), Google Cloud Platform or GCP (<https://cloud.google.com>), and Oracle Cloud Infrastructure or OCI (<https://www.oracle.com/cloud>). The units of

computational power offered by CSPs receive different names such as instance, vm or even server. The expression 'instance' has become the most widely used by cloud practitioners and hereafter refers to a minimum of computational power, storage, and a network connection. CSPs offer instances in many sizes that can be as small as a shared core. However, the interest here resides on the largest instances that use all the cores available in the socket(s). This is the most common situation for WRF as well as for many HPC apps. Table 1 lists the characteristics of several processors used by these CSPs in their series targeting computational power. In addition to maker, model and architecture, Table 1 indicates whether the instances powered by these processors are single or dual sockets and the number of cores per socket.

TABLE 1- Custom processors for cloud environments

Maker & model	Arch.	Instance socket	Cores per socket	CSP
AWS Graviton2	Aarch64	Single	64	AWS
AMD EPYC™ 7R32	x86_64	Single	48	AWS
AMD EPYC™ 7V12	x86_64	Dual	56	GCP
Ampere® Altra™	Aarch64	Dual	80	Available Q4/2020
Intel® Xeon® Platinum 8124	x86_64	Dual	18	AWS
Intel® Xeon® Scalable (GCP)	x86_64	Dual	15	GCP

The results discussed here measure system performance using the WRF component of the NWSC-3 benchmark suite. NCAR has developed this benchmark, which is based on version 4.1 of WRF, to assess the next generation of High-Performance Computing and Storage System. The main objective of this new benchmark is to facilitate performance comparison between systems with different hardware or even based on different architectures. To accomplish this objective, NCAR decreased the size of data sets

versus previous benchmark version and stressed the measurement of scalability for systems powered with many cores. Fig. 1 illustrates the performance of the processors listed in Table 1 in a dual-socket configuration. The final figure represents the time needed to complete the benchmark. AWS offers Graviton2 and EPYC™ 7R32 as single-socket instances so the benchmarks for these processors have used a two instance configuration.

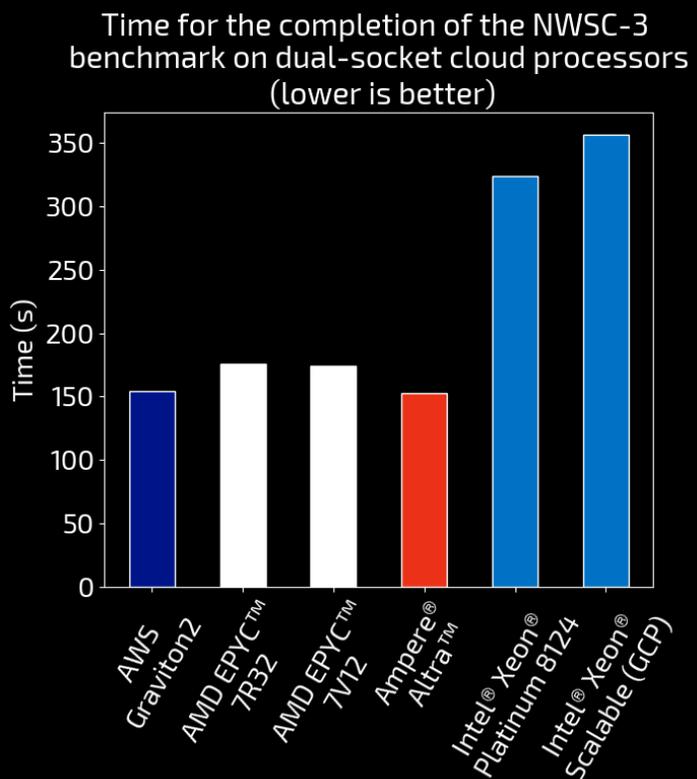


Fig. 1 - Comparative performance of WRF on dual-socket cloud processors.

### Performance of WRF on IaaS clusters

The performance of large-scale simulations requires clusters composed by hundreds or even thousands of cores connected through a high bandwidth network. The assessment of WRF on IaaS clusters involve the deployment of several instances to reach the desired core count. To have a reference for on-premises performance, the present assessment uses

benchmarks performed at NCAR's Cheyenne supercomputer. Cheyenne occupies the 52nd position in the list ranking the top systems worldwide ([www.top500.org](http://www.top500.org)) as of June of 2020. The first level of the benchmarks performed at Cheyenne uses 576 cores and the IaaS clusters have been built around this figure. Table 2 lists the total number of cores and other cluster characteristics.

TABLE 2 - Instance characteristics of IaaS clusters for WRF benchmarking

CSP	Instance	Processor
AWS	c5.18xlarge	Intel® Xeon® Platinum 8124
AWS	c5n.18xlarge	Intel® Xeon® Platinum 8124
AWS	c5a.24xlarge	AMD EPYC™ 7R32
AWS	c6g.16xlarge	AWS Graviton2
GCP	c2-standard-60	Intel® Xeon® Scalable (GCP)
GCP	n2d-standard-60	AMD EPYC™ 7V12
OCI	BM.HPC2.36	Intel® Xeon® Gold 6154

In addition to pure computational power, WRF performance depends on the network bandwidth and latency connecting the instances or nodes. In the case of IaaS clusters, CSPs provide the ability to create a subnet within a private network interconnecting

cluster instances. The hardware connections to the instances can be either Ethernet or Infiniband (IB) switches similarly to on-premises systems, which leads to two communication standards. Additionally, some CSPs are introducing their own proprietary technologies to network communication. An example is AWS's Elastic Fabric Adapter (EFA) that works as a network interface to decrease latencies when many instances are interconnected. Table 3 provides a summary of the network properties used to build the IaaS clusters.

TABLE 3 - Network properties interconnecting the IaaS clusters used for WRF benchmarking

Cluster	Network type	Maximum bandwidth
AWS - c5.24xlarge	Ethernet	25 Gbps
AWS - c5n.18xlarge	EFA	100 Gbps
AWS - c5a.24xlarge	Ethernet	20 Gbps
AWS - c6g.16xlarge	Ethernet	25 Gbps
GCP - c2.standard.60	Ethernet	32 Gbps
GCP - n2d.standard.224	Ethernet	32 Gbps
OCI - BM.HPC2.36	IB	100 Gbps

Figure 2 shows the total computational time for the NWSC-3 benchmark reported for Cheyenne

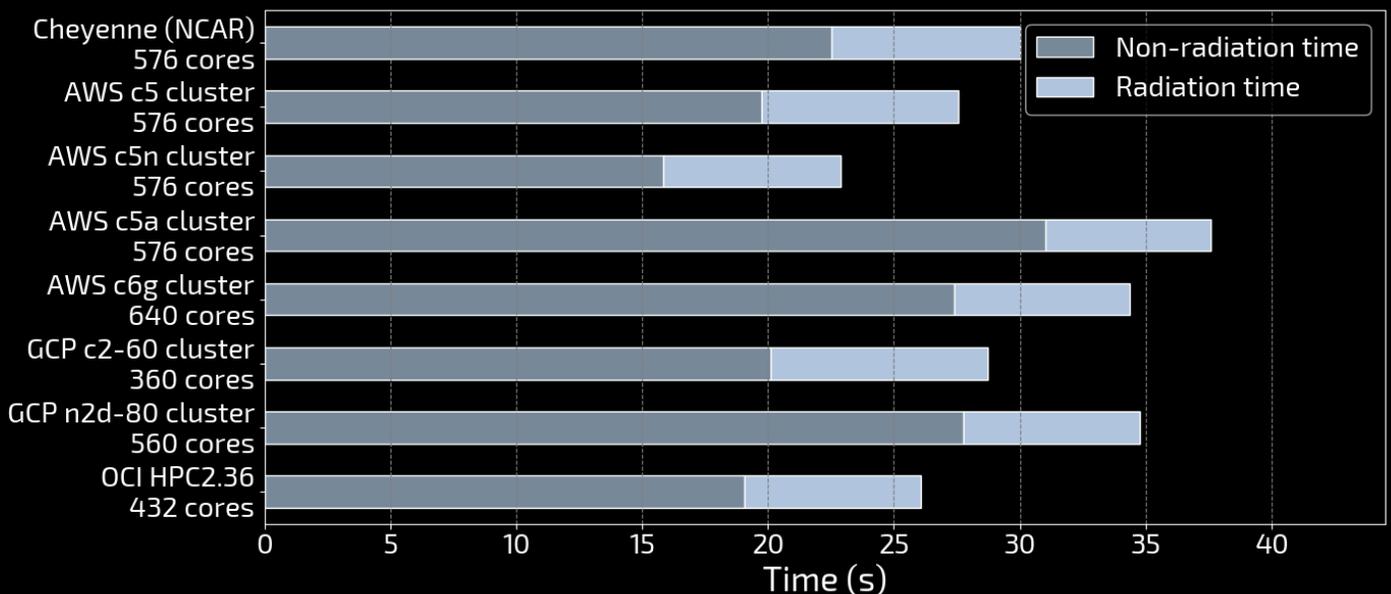


Figure 2 - WRF performance (NWSC-3 benchmark) on IaaS clusters from the public cloud.

as well as the measured times for each IaaS cluster. The total computational time is composed by radiation time (~25%) plus non-radiation time (~75%). The times achieved with IaaS clusters are competitive versus on-premises benchmarks. Furthermore, some of the IaaS clusters with newer processors outperform the Cheyenne benchmarks.

### Cost

In addition to performance, any decision to migrate WRF to IaaS from the public cloud must include cost considerations. We strongly recommend a total cost of ownership (TCO) evaluation for organizations interested in this migration as several factors dictate the final cost. Cost estimates must contemplate that CSPs provide a wide variety of payment options. To perform an initial estimate, the cost analysis considers the following 3 basic tiers:

- On-demand price: This is full price offered by CSPs that do not require any

commitment and offer complete flexibility.

- Reserved price: This price requires a commitment for a specific period. The discounts depend on many factors such as the type of resources, terms of commitment or when the payments are made. The present evaluation considers a standard one-year commitment.
- Spot price: This mode provides the maximum discount, but instances can be preempted by the CSP.

Fig. 3 provides a cost comparison for several clusters. The cost estimate only accounts for computational power and has been normalized based on the ratio between price and performance. In addition to any potential savings from using reserved or spot instances, Fig. 3 also shows that Aarch64 IaaS is not only closing the performance gap versus x86\_64 hardware, but it also has the potential to become more cost-effective in the near future.

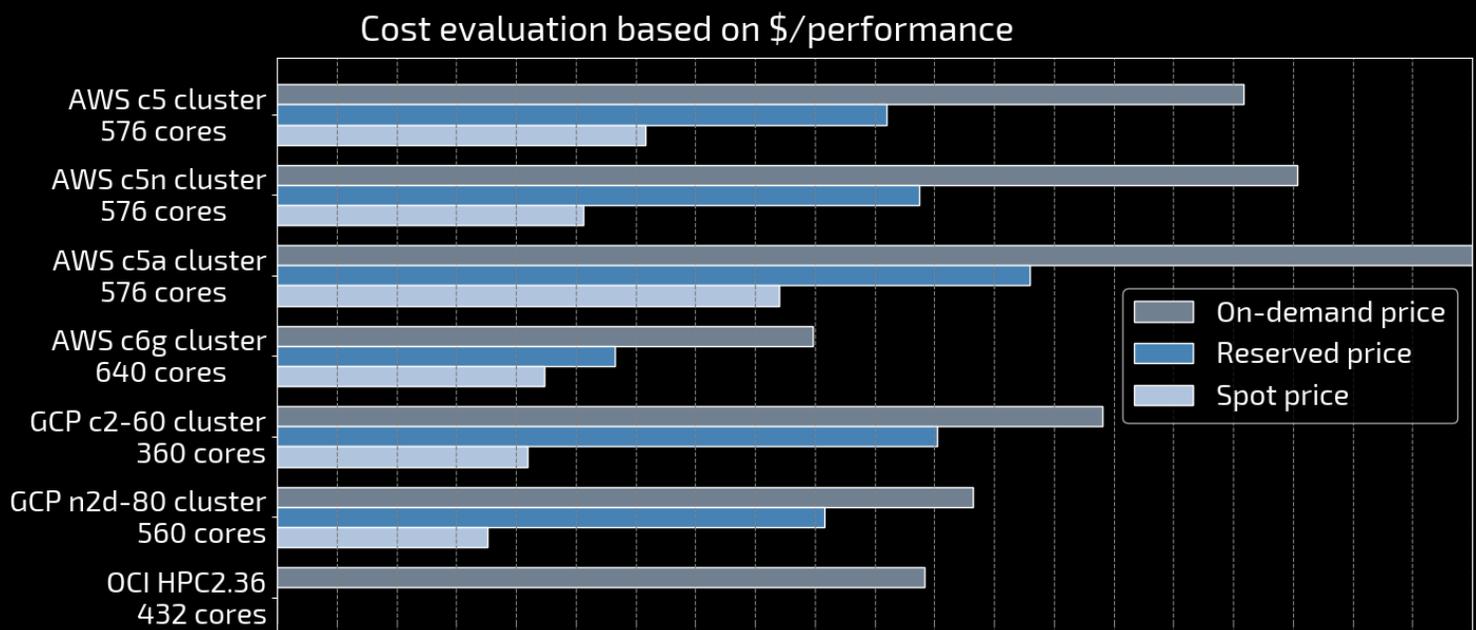


Figure 3 - Cost comparison for WRF on IaaS clusters from the public cloud.

## CONCLUSIONS

Evaluation of WRF in the public cloud has shown new IaaS capabilities to be able to achieve performance levels traditionally associated with supercomputers or clusters. The flexibility offered by IaaS, which facilitates the configuration of hardware from a few cores to cloud clusters composed by hundreds or even thousands of cores, expedites tailored configuration for different WRF predictions. Furthermore, the introduction of new processors with a high core count and specifically developed for cloud environments, such as Ampere® Altra™ or AWS Graviton2, results in better per instance or socket performance, and the potential for economic savings. Lastly, an element not to be overlooked by WRF users is the practically unlimited storage capacity offered by CSPs, and that balances the needs for computational power and large storage capabilities.

---

## Contact information

Benchmarks performed by odyhpc.com. For questions or more information, contact us at:

ODYHPC.COM - Odysseus Computational Solutions

Murrysville, PA 15668

USA

TEL: +1- 724-647-7367

e-mail: support@odyhpc.com



---

### DISCLAIMER

The information contained herein is for informational purposes only and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions, and typographical errors. ODYHPC.COM makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of hardware, software or other products described herein.