

Standing up an HPC cluster on AWS using AWS ParallelCluster

Angel Pizarro, HPC @ AWS MUG '21



The metric for success should be time-to-results



Massive capacity when needed to speed up time to Finite capacity, usually with long queues to results, and agile environment when additional wait in. hardware and software experimentation is needed.

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

Source: Hyperion ROI Study (http://www.hyperionresearch.com/roi-with-hpc/)







HPC workloads across industries



Life Sciences



Financial Services





Design & Engineering







© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.



Autonomous Vehicles



HPC workloads with different compute and throughput characteristics





What do HPC customers need

High Performance

Reliability

Low Latency/Overhead

Consistency & Jitter free network

Fairness



Scalable Reliable Datagram (SRD)

A reliable high-performance lower-latency network transport

Guaranteed delivery

Does not consume any resources on EC2 instances ٠

Network aware multipath routing

Optimally utilizes all network paths (ECMP), no hot-spots •

Orders of magnitude lower tail latency & jitter

Fast recovery from network events ٠

No ordering guarantees

No head-of-line blocking ٠

https://ieeexplore.ieee.org/document/9167399







Elastic Fabric Adapter – Networks built to scale





- OS bypass \checkmark
- **GPUdirect and RDMA** \checkmark
- Libfabric core supports \checkmark wide array of MPIs and NCCL



OCMP-enabled packet spraying and cloud-scale congestion control

ANSYS Fluent External flow over F1 race car (140M cell mesh)



scaling efficiency vs ~48% using C5 w/o EFA

C5a

At ~3,000 cores (~83 nodes), C5n+EFA shows ~89%



HPC software stack on Amazon EC2





Elastic Fabric Adapter (EFA)

Scale tightly coupled HPC applications on AWS



EFA AWS HPC/ML Network Interface Instance flexibility Infrastructure elasticity High data throughput Low-latency message passing Faster application time-tocompletion



High Performance Computing (HPC) on AWS

Virtual Private Cloud on AWS

A

 (\mathbf{A})

3D graphics virtual workstation

License managers and cluster head

nodes with job schedulers

On AWS, secure and welloptimized HPC clusters can be automatically created, operated, and torn down in just minutes





© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.



Thin or zero client no local data







Integrated with AWS services you need



© 2021, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark

One-stop shop to set up your HPC cluster



AWS Batch



Running EFA using AWS ParallelCluster

Subtitle



Getting started using Spack with AWS ParallelCluster



G Github

https://workshops.aws/categories/HPC

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Traden

up a AWS ParallelCluster cluster with C6g and C6gn instances with Spack, and ReFrame preinstalled.

Info!

workshop).

Target Audience

benchmark code on Graviton2 instances.

Background Knowledge

You should be familiar with HPC setups in general (specifically SLURM), have some knowledge about the linux command line to walk through the example.

To make the most out of it and get your hands dirty you will need to learn or know about Spack and python to optimize codes.

Duration



AWS ParallelCluster architecture



DEMO TIME (sort of 🥞)



pizarroa@147ddacf1e59:~/src/mug21-pcluster

- → mug21-pcluster asciinema play -s 4 pcl-conf.cast
- → mug21-pcluster p

. . .

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

T#1



Feature: Multiple Slurm Queues



12	[scaling demo]	
13	<pre>scaledown_idletime = 5</pre>	#
14		
15	[cluster multi-queue]	
16	<pre>key_name = pc-key-mug21</pre>	
17	<pre>base_os = alinux2</pre>	
18	scheduler = slurm	
19	<pre>master_instance_type = c5.2xlarge</pre>	
20	<pre>vpc_settings = default</pre>	
21	<pre>scaling_settings = demo</pre>	
22	<pre>queue_settings = spot,ondemand</pre>	
23		
24	[queue spot]	
25	<pre>compute_resource_settings = spot_i1;</pre>	, sp
26	<pre>compute_type = spot</pre>	#
27		
28	[compute_resource spot_i1]	
29	<pre>instance_type = c5.xlarge</pre>	
30	<pre>min_count = 0</pre>	#
31	<pre>max_count = 10</pre>	#
32		
-33	[compute_resource spot_i2]	
34	<pre>instance_type = t3.medium</pre>	
35	<pre>min_count = 1</pre>	
36	<pre>initial_count = 2</pre>	
37		
38	[queue ondemand]	
39	<pre>compute_resource_settings = ondemand</pre>	d_i
40	<pre>disable_hyperthreading = true</pre>	#

https://docs.aws.amazon.com/parallelcluster/latest/ug/tutorial-mqm.html

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

optional, defaults to 10 minutes

ot_i2
optional, defaults to ondemand

optional, defaults to 0 optional, defaults to 10

1 optional, defaults to false



Feature: Amazon FSx for Lustre

Parallel file system



100+ GiB/s throughput Millions of IOPS Consistent sub-millisecond latencies

SSD-based





Supports hundreds of thousands of cores

T -1	
15	[cluster default
16	key_name = pc-ke
17	scheduler = slur
18	<pre>master_instance_</pre>
19	<pre>base_os = alinux</pre>
20	<pre>vpc_settings = d</pre>
21	<pre>queue_settings =</pre>
22	scaling settings
23	<pre>fsx_settings = f</pre>
24	
25	[fsx fs-mug21]
26	<pre>shared_dir = /fs</pre>
27	<pre>storage_capacity</pre>
28	imported_file_ch
29	<pre>export_path = s3</pre>
30	<pre>import_path = s3</pre>
31	weekly_maintenan

22

:] ey-mug21 m type = c5.2xlarge 2 lefault : compute = demo s-mug21

х

= 3600

unk_size = 1024 ://bucket/folder ://bucket ce_start_time = 1:00:00



EC2 UltraClusters of P4d

Supercomputing-class performance for deep learning workflows



© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

- to large clusters for distributed training
- lower cost to train
- EC2 UltraClusters with to over 4,000 GPUs

• Based on NVIDIA A100 GPUs

Availability to scale-out

 2.5x better deep learning performance and 60%

EFA enables 400 Gbps and allows you to scale



Thank you!

https://docs.aws.amazon.com/parallelcluster/latest/ug/what-is-awsparallelcluster.html

> https://workshops.aws/categories/HPC https://spack-tutorial.workshop.aws/ https://www.hpcworkshops.com/

https://mvapich.cse.ohio-state.edu/userguide/mv2x-aws/ https://mvapich.cse.ohio-state.edu/userguide/userguide_spack/



