



Achieving High Cluster Utilization with Over-Quota GPU Allocation

NetApp Solutions

Kevin Hoke
May 24, 2021

This PDF was generated from https://docs.netapp.com/us-en/netapp-solutions/ai/osrunai_achieving_high_cluster_utilization_with_over-uota_gpu_allocation.html on October 21, 2021. Always check docs.netapp.com for the latest.

Table of Contents

Achieving High Cluster Utilization with Over-Quota GPU Allocation 1

Achieving High Cluster Utilization with Over-Quota GPU Allocation

In this section and in the sections [Basic Resource Allocation Fairness](#), and [Over-Quota Fairness](#), we have devised advanced testing scenarios to demonstrate the Run:AI orchestration capabilities for complex workload management, automatic preemptive scheduling, and over-quota GPU provisioning. We did this to achieve high cluster-resource usage and optimize enterprise-level data science team productivity in an ONTAP AI environment.

For these three sections, set the following projects and quotas:

Project	Quota
team-a	4
team-b	2
team-c	2
team-d	8

In addition, we use the following containers for these three sections:

- Jupyter Notebook: `jupyter/base-notebook`
- Run:AI quickstart: `gcr.io/run-ai-demo/quickstart`

We set the following goals for this test scenario:

- Show the simplicity of resource provisioning and how resources are abstracted from users
- Show how users can easily provision fractions of a GPU and integer number of GPUs
- Show how the system eliminates compute bottlenecks by allowing teams or users to go over their resource quota if there are free GPUs in the cluster
- Show how data pipeline bottlenecks are eliminated by using the NetApp solution when running compute-intensive jobs, such as the NetApp container
- Show how multiple types of containers are running using the system
 - Jupyter Notebook
 - Run:AI container
- Show high utilization when the cluster is full

For details on the actual command sequence executed during the testing, see [Testing Details for Section 4.8](#).

When all 13 workloads are submitted, you can see a list of container names and GPUs allocated, as shown in the following figure. We have seven training and six interactive jobs, simulating four data science teams, each with their own models running or in development. For interactive jobs, individual developers are using Jupyter Notebooks to write or debug their code. Thus, it is suitable to provision GPU fractions without using too many cluster resources.

```

root@run-deploy:~# kubectl get pods -A
NAME                                STATUS    AGE      NODE              IMAGE                                     TYPE      PROJECT  USER  GPUS  CREATED BY  CLI  SERVICE URL(S)
b-4-gg                             Running   2m       dgx1-2            gcr.io/run-ai-demo/quickstart          Train     team-b   root   2     true      team-b     true
c-5-g                               Running   2m       dgx1-2            gcr.io/run-ai-demo/quickstart          Train     team-c   root   1     true      team-c     true
c-4-gg                             Running   2m       dgx1-1            gcr.io/run-ai-demo/quickstart          Train     team-c   root   2     true      team-c     true
b-3-g                               Running   2m       dgx1-1            gcr.io/run-ai-demo/quickstart          Train     team-b   root   1     true      team-b     true
c-3-g02                            Running   2m       dgx1-1            gcr.io/run-ai-demo/quickstart          Interactive team-c   root   0.2   true      team-c     true
d-1-gggg                           Running   2m       dgx1-2            gcr.io/run-ai-demo/quickstart          Train     team-d   root   4     true      team-d     true
c-2-g03                            Running   2m       dgx1-1            gcr.io/run-ai-demo/quickstart          Interactive team-c   root   0.3   true      team-c     true
c-1-g05                            Running   2m       dgx1-1            gcr.io/run-ai-demo/quickstart          Interactive team-c   root   0.5   true      team-c     true
a-2-gg                             Running   3m       dgx1-1            gcr.io/run-ai-demo/quickstart          Train     team-a   root   2     true      team-a     true
b-2-g04                            Running   3m       dgx1-2            gcr.io/run-ai-demo/quickstart          Interactive team-b   root   0.4   true      team-b     true
a-1-g                               Running   3m       dgx1-1            gcr.io/run-ai-demo/quickstart          Train     team-a   root   1     true      team-a     true
b-1-g06                            Running   3m       dgx1-2            gcr.io/run-ai-demo/quickstart          Interactive team-b   root   0.6   true      team-b     true
a-1-1-jupyter                      Running   3m       dgx1-1            jupyter/base-notebook                  Interactive team-a   root   1     true      team-a     true
https://10.61.218.134/a-1-1-jupyter

```

The results of this testing scenario show the following:

- The cluster should be full: 16/16 GPUs are used.
- High cluster utilization.
- More experiments than GPUs due to fractional allocation.
- team-d is not using all their quota; therefore, team-b and team-c can use additional GPUs for their experiments, leading to faster time to innovation.

Next: Basic Resource Allocation Fairness

Copyright Information

Copyright © 2021 NetApp, Inc. All rights reserved. Printed in the U.S. No part of this document covered by copyright may be reproduced in any form or by any means-graphic, electronic, or mechanical, including photocopying, recording, taping, or storage in an electronic retrieval system-without prior written permission of the copyright owner.

Software derived from copyrighted NetApp material is subject to the following license and disclaimer:

THIS SOFTWARE IS PROVIDED BY NETAPP "AS IS" AND WITHOUT ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE, WHICH ARE HEREBY DISCLAIMED. IN NO EVENT SHALL NETAPP BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

NetApp reserves the right to change any products described herein at any time, and without notice. NetApp assumes no responsibility or liability arising from the use of products described herein, except as expressly agreed to in writing by NetApp. The use or purchase of this product does not convey a license under any patent rights, trademark rights, or any other intellectual property rights of NetApp.

The product described in this manual may be protected by one or more U.S. patents, foreign patents, or pending applications.

RESTRICTED RIGHTS LEGEND: Use, duplication, or disclosure by the government is subject to restrictions as set forth in subparagraph (c)(1)(ii) of the Rights in Technical Data and Computer Software clause at DFARS 252.277-7103 (October 1988) and FAR 52-227-19 (June 1987).

Trademark Information

NETAPP, the NETAPP logo, and the marks listed at <http://www.netapp.com/TM> are trademarks of NetApp, Inc. Other company and product names may be trademarks of their respective owners.