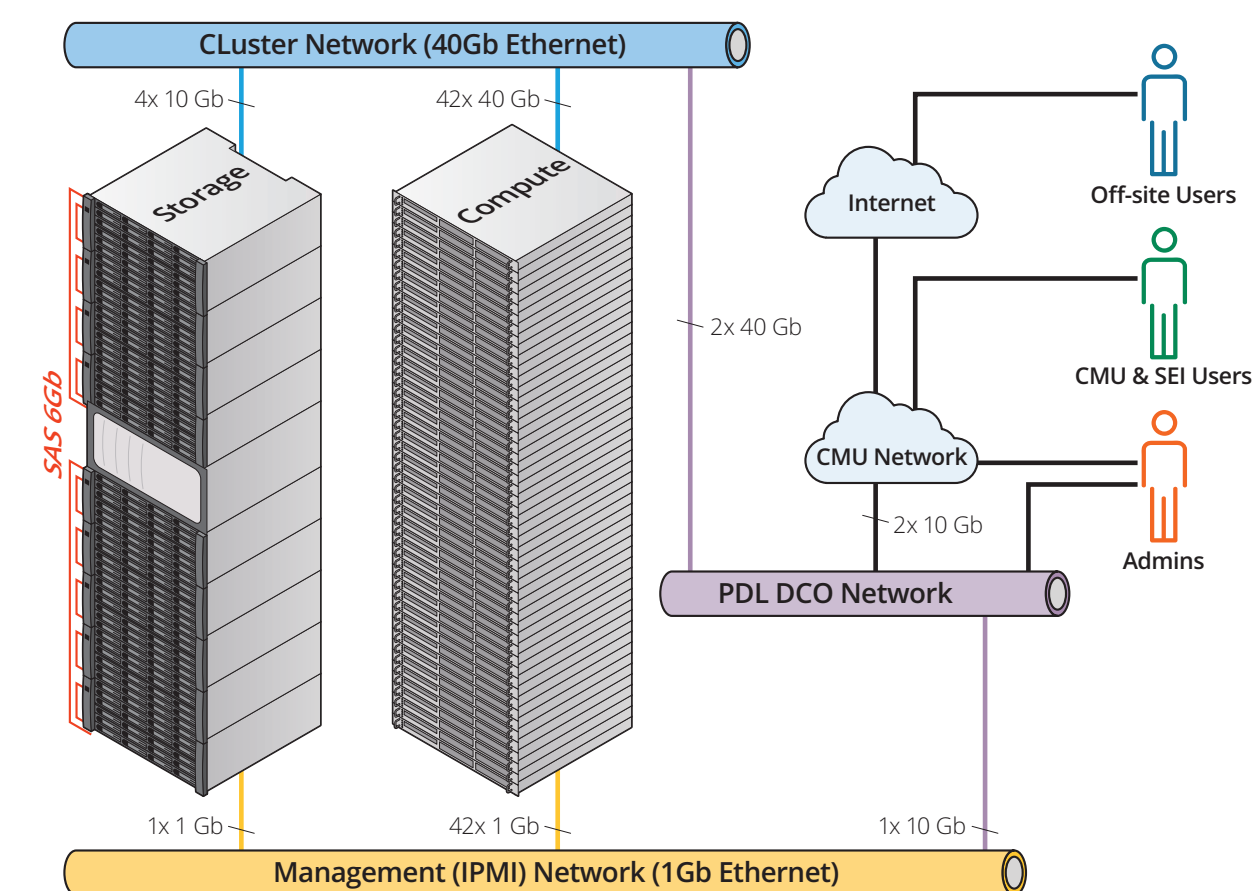


Measuring Performance of Big Learning Workloads

Big Learning platforms—large scale hardware and software systems designed to perform large-scale machine learning (ML) workloads on big data—are extremely complex and lack consistent and sufficient reporting of performance metrics. These difficulties can slow DoD adoption of new advances in machine learning. A key obstacle to overcome is in the collection and analysis of these metrics.



Component	Per Node	Totals (42 nodes)
CPU, Xeon E5, 2 GHz	16 cores, 32 threads	672 cores, 1344 threads
GPU, Titan X	3072 cores, 12 GB RAM	129K cores
RAM, DDR4	64 GB	2688 GB
NVMe storage	400 GB	16.8 TB
HDD storage	8 TB	336 TB
Network	40 GB	
Persistent storage		432 TB

The Big Learning cluster consists of 42 compute nodes each with CPU and GPU processing units, complex storage system and fast networking. It supports research in the development of parallel ML computing frameworks as well as development and testing of large-scale metrics collection systems.

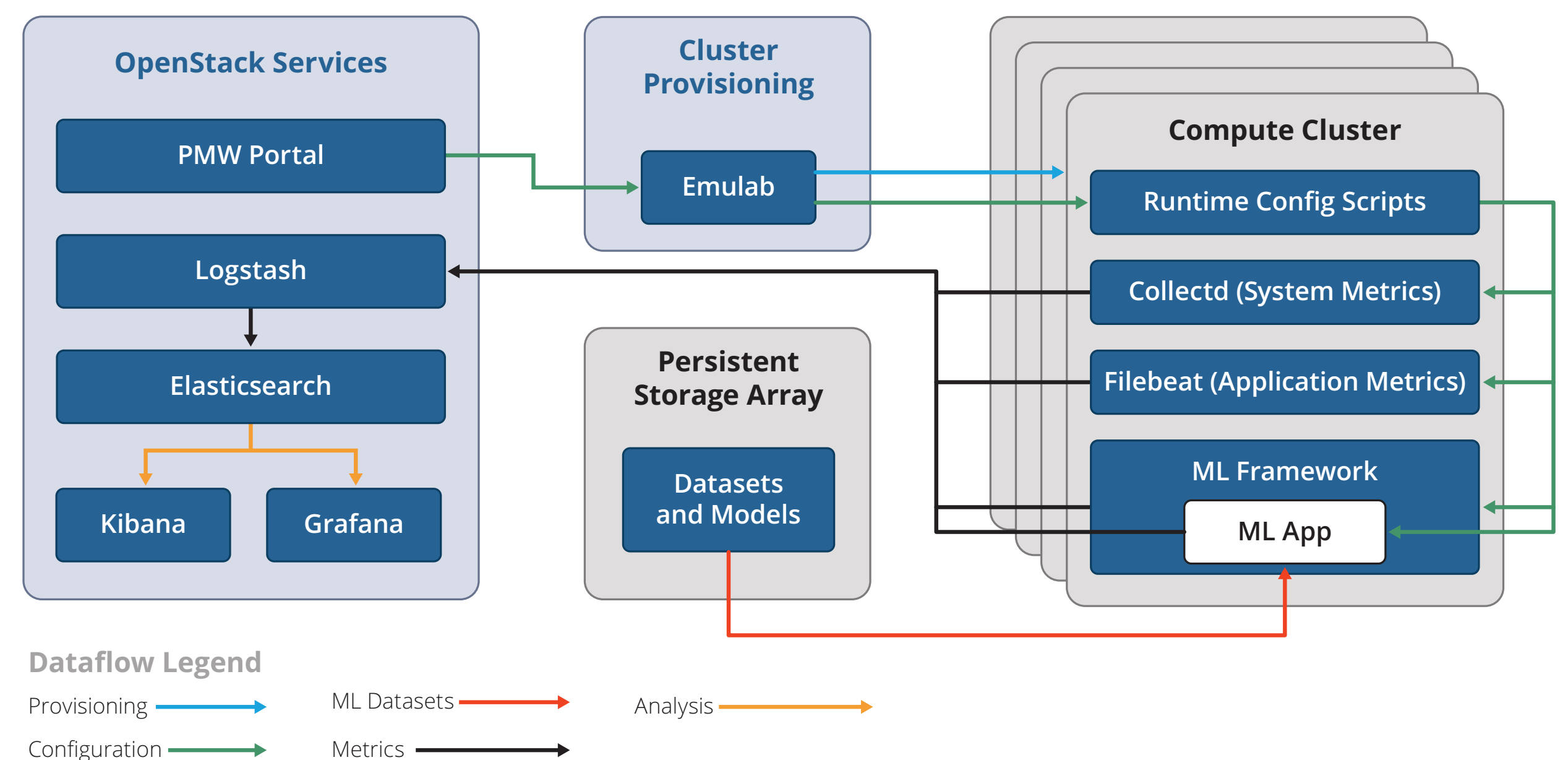
Big Learning Cluster. The Big Learning cluster, located at the Carnegie Mellon University Parallel Data Lab, was designed to support research in evaluating existing Big Learning platforms and developing new platforms for a wide variety of large-scale ML applications. It is a distributed cluster with CPU and GPU processors, a complex storage hierarchy, a high-bandwidth/low latency network for communication, and a large persistent store.

Performance Measurement Workbench. We developed the Performance Measurement Workbench (PMW) to collect metrics about the performance of hardware components (CPU, memory, disk, network, etc.) and the performance of the ML algorithms (accuracy, convergence, iteration times, etc.) that run on the Big Learning cluster.

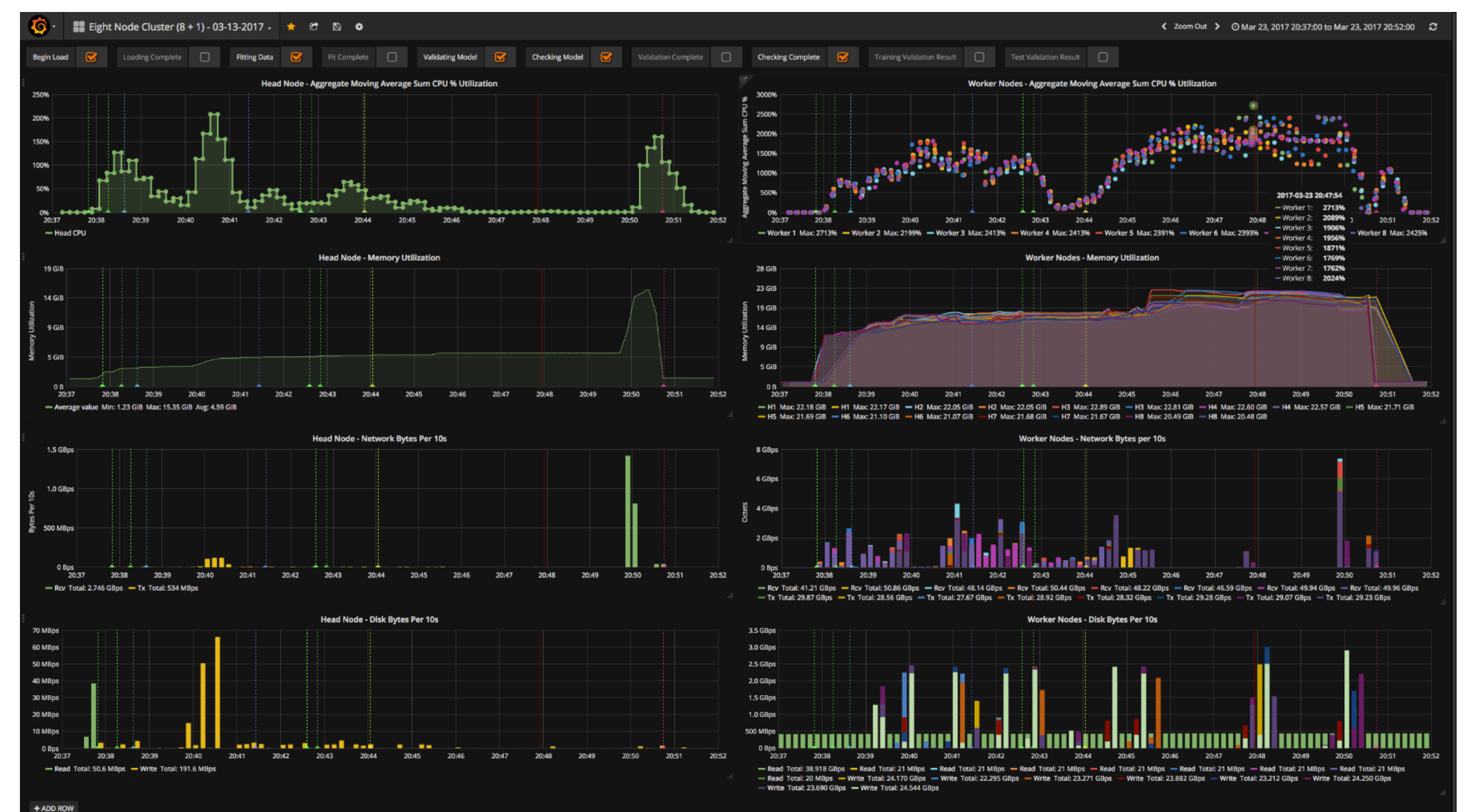
Goal: Ease of Use. PMW provides a simple, web-based portal for researchers and users to configure and submit jobs, collect and store metrics, and analyze data both during computation and in post mortem.

Goal: Reproducibility. PMW not only collects the performance metrics for each job, but it can also collect and store the configuration of the operating system, the ML platform, and the algorithm being run. With this information reproducible experiments are achievable.

With the Performance Measurement Workbench—combining a few open-source software packages (especially Elastic Stack)—we have demonstrated how consistent and complete measurement metrics for complex Big Learning systems can be collected. PMW has the added benefit of supporting collection of configuration aimed at reproducibility of results.



Performance Measurement Workbench system architecture. PMW provides a simple, web-based portal for submitting jobs, operating system images with collection tools preconfigured, and persistent database query and visualization services using the open-source Elastic Stack.



PMW's dashboard display using Grafana integrates with Elastic Stack to achieve complex visualizations. This example shows system metrics for a Spark MLlib job that uses one "head" node (displayed on the left) and eight worker nodes to perform a logistic regression algorithm (displayed on the right).

Copyright 2017 Carnegie Mellon University. All Rights Reserved.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

Internal use:* Permission to reproduce this material and to prepare derivative works from this material for internal use is granted, provided the copyright and "No Warranty" statements are included with all reproductions and derivative works.

External use:* This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other external and/or commercial use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

* These restrictions do not apply to U.S. government entities.

Carnegie Mellon® is registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM17-0733
Measuring Performance of Big Learning Workloads