# HADOOP INFRASTRUCTURE SCALING WITH THE DELL POWEREDGE FX2
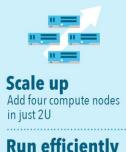


**BIG DATA, SMALL FOOTPRINT**

Hadoop® on Dell™ PowerEdge™ FX2

Get **2.25X** the performance by adding a second FX2

**Scale up**
Add four compute nodes in just 2U

**Run efficiently**
with a balanced use of hardware resources

**Get results faster**
Finish workloads in less than half the time

Dell PowerEdge FC430 server

Dell PowerEdge FD332 storage node

Powered by the Intel® Xeon® processor E5-2600 v3 product family

When wading into the Hadoop big data pool, it's important to select a solution that can handle the jobs you run, and one that is flexible enough to scale well as the size of your big data needs increase over time. The Dell PowerEdge FX2 is a datacenter solution that combines all the essential IT elements—servers, storage, and networking blocks—into a very compact 2U chassis. You can tailor the Dell PowerEdge FX2 solution to meet your unique workload needs, such as Hadoop workloads that process big data. In particular, Hadoop thrives with uniform compute scale-out and a high disk-to-compute ratio for Hadoop File System (HDFS) storage capacity, both of which the Dell PowerEdge FX2 provides.

In the Principled Technologies labs, we tested a single Dell PowerEdge FX2 with four PowerEdge FC430 nodes, and found that it completed our Hadoop workload in 25 minutes and 58 seconds. When we added a second Dell PowerEdge FX2, Hadoop performance scaled well: by just adding a second FX2 cluster, it cut the job time by more than half. All the way down to 11 minutes and 31 seconds.

While many Hadoop infrastructures have dozens of nodes, you want to be sure when starting out to choose a flexible and scalable solution. By choosing the Dell PowerEdge FX2 to start your Hadoop infrastructure, you can get all the benefits of its unique converged infrastructure design, which can include fast performance, simplified management, and space savings thanks to its dense nature. And when you decide it's time to scale out your solution, adding a cluster and cutting job times in half is simple thanks to the Dell PowerEdge FX2 all-in-one chassis.

# BIG DATA IN SMALL SPACES

Sorting and reorganizing the data you collect can help your organization get a handle on how your business runs. Hadoop is an application that breaks big data into smaller sets and spreads them out over multiple server nodes, making big data analysis fast and scalable.

The Dell PowerEdge FX2 solution configured with four server nodes and two storage blocks can run Hadoop workloads, and does it all in just 2U of space. With servers, storage, and networking sharing a common chassis, the Dell PowerEdge FX2 brings all the elements of a traditional datacenter into a single chassis, which can simplify your infrastructure. Because the PowerEdge FX2 can support a number of different configurations of those elements, you can build your organization's PowerEdge FX2 to fit your exact workload needs. These are just some of the kinds of benefits that the Dell PowerEdge FX2 can bring to organizations that traditional server and storage setups can't; it helps you make the most efficient use of each element in your infrastructure.

# WHAT WE FOUND
## About the results

Our test workload used 300GB of data and performed several common Hadoop operations on large datasets, including data generation, sorting the data, and data validation. Our workload executed a short data integrity check after the data generation and sorting portions. These operations are simple but highly representative of real-world Hadoop workloads that stress the Map-Reduce framework and the Hadoop Filesystem API.

We used Cloudera Distributed Hadoop (CDH) 5.4.2 as our Hadoop cluster software. We set up the first Dell PowerEdge FX2 to house the Edge, Name, and Data Node roles across four nodes. The second Dell PowerEdge FX2 unit had four Data Nodes.[1]

We tested the scalability of the Dell PowerEdge FX2 with four Dell PowerEdge FC430 nodes and two Dell PowerEdge FD332 storage arrays by running the TPCx-HS 300GB workload on one Dell PowerEdge FX2, then adding a second Dell PowerEdge FX2 with the same hardware configuration and measuring the time required to run the same workload. When we added a second Dell PowerEdge FX2 to the cluster, the workload time decreased by 56 percent (see Figure 1).

---

[1] For specific Hadoop tuning parameters, see our full report at
http://www.principledtechnologies.com/Dell/FX2_Hadoop_scaling_1015.pdf.

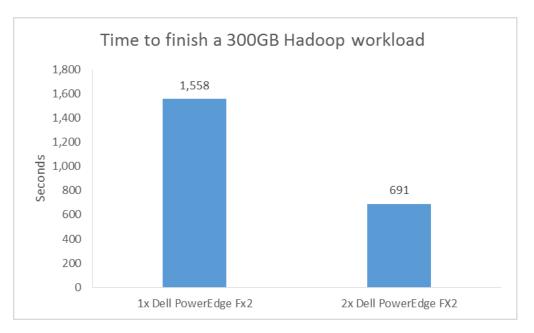## Time to finish a 300GB Hadoop workload



**Figure 1: Time to complete our Hadoop workload, in seconds.**

## Efficient use of resources

A properly tuned Hadoop cluster can take advantage of all the hardware subsystems (CPU, memory, and storage) you make available to it. Based on Hadoop example workloads TeraGen, TeraSort, and TeraValidate, our workload was dependent on CPU, memory and disk resources, so it was important that all three subsystems were adequately utilized.

Not only did the Dell PowerEdge FX2 unit show excellent scaling, it was also able to provide balanced use of its hardware resources in both phases of testing. Because each of the balanced utilization, an owner of a similarly configured Dell PowerEdge FX2 could run this workload confident that resources are being used efficiently. That same owner could then purchase a second, identical Dell PowerEdge FX2 and be comfortable knowing that their workloads continue to operate without leaving idle hardware on the table.

Figure 2 through 4 show the utilization metrics (averaged across the Data Nodes for each phase) of each hardware subsystem during the first and second phases of our testing.

As Figure 2 shows, CPU utilization remained high for every portion of the workload during the first phase of testing. Adding a second Dell PowerEdge FX2 did not change the CPU utilization performance profile, showing that this workload scales well from a CPU perspective. The slight decrease in CPU activity during sorting is due to the disk-intensive reduce portion of that operation.
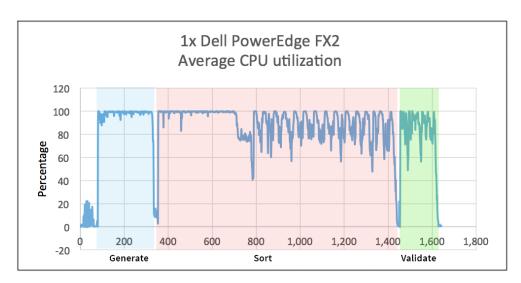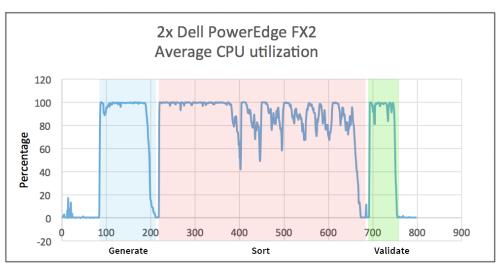
**Figure 2: Average Data Node CPU utilization percentages for 1x Dell PowerEdge FX2 and for 2x Dell PowerEdge FX2.**

We tuned our Hadoop cluster to take full advantage of the available memory in each node. As Figure 3 shows, the workload was able to make use of all the memory in both phases of testing, indicating that the workload scales well from a memory usage perspective.
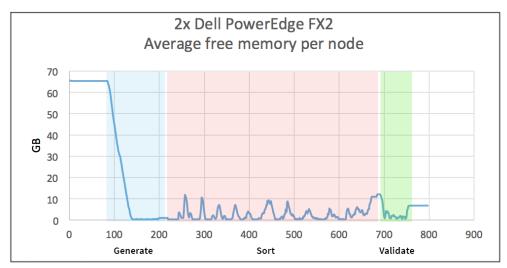


**Figure 3: Average free memory per Data Node for 1x Dell PowerEdge FX2 and for 2x Dell PowerEdge FX2.**

Disk performance is critical to many Hadoop operations, and the three major operations in our workload are no exception. The Dell PowerEdge FD332 storage blocks and shared RAID controllers allow presentation of the disks in RAID or HBA mode. While a RAID group can add performance and data replication for many common workloads, Hadoop prefers HBA mode as the Hadoop Distributed File System (HDFS) handles replication. Our workload was able to fully utilize the disks during data generation and the reduce portion of the sorting operations. These operations occur in memory whenever possible, which means that disk utilization decreases during data validation

and the map portion of sorting As Figure 4 shows, the level of disk utilization was similar in both phases of testing, indicating good scaling of disk resources.
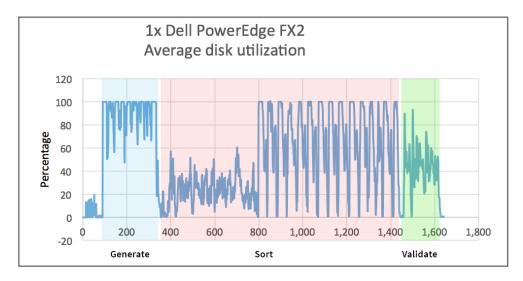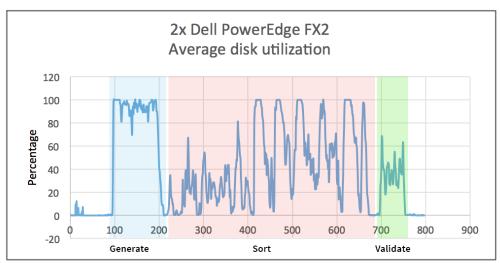


**Figure 4: Average disk utilization across all Data Nodes for 1x Dell PowerEdge FX2 and for 2x Dell PowerEdge FX2.**

# CONCLUSION

The definition of a successful Hadoop solution need not be limited to whether or not the hardware can run the jobs and sort the data. As our tests show, the Dell PowerEdge FX2 was powerful enough to run our Hadoop workload, but more importantly, it scaled well when we added another cluster. Adding a second PowerEdge FX2 chassis complete with four Dell PowerEdge FC430 server nodes and Dell PowerEdge FD332 storage cut the time to run our Hadoop job in half. The all-in-one chassis that brings compute, storage, and networking together can also offer other benefits inherent in its design: the Dell PowerEdge FX2 can sort big data in a small space, which can also deliver space savings and ease the burden of managing the Hadoop solution.

# ABOUT PRINCIPLED TECHNOLOGIES

**Principled Technologies®**

Principled Technologies, Inc.
1007 Slater Road, Suite 300
Durham, NC, 27703
www.principledtechnologies.com

We provide industry-leading technology assessment and fact-based marketing services. We bring to every assignment extensive experience with and expertise in all aspects of technology testing and analysis, from researching new technologies, to developing new methodologies, to testing with existing and new tools.

When the assessment is complete, we know how to present the results to a broad range of target audiences. We provide our clients with the materials they need, from market-focused data to use in their own collateral to custom sales aids, such as test reports, performance assessments, and white papers. Every document reflects the results of our trusted independent analysis.

We provide customized services that focus on our clients' individual requirements. Whether the technology involves hardware, software, Web sites, or services, we offer the experience, expertise, and tools to help our clients assess how it will fare against its competition, its performance, its market readiness, and its quality and reliability.

Our founders, Mark L. Van Name and Bill Catchings, have worked together in technology assessment for over 20 years. As journalists, they published over a thousand articles on a wide array of technology subjects. They created and led the Ziff-Davis Benchmark Operation, which developed such industry-standard benchmarks as Ziff Davis Media's Winstone and WebBench. They founded and led eTesting Labs, and after the acquisition of that company by Lionbridge Technologies were the head and CTO of VeriTest.