

HPC Applications Performance on R740 with V100 GPUs

Authors: Frank Han, Rengan Xu, Nishanth Dandapanthula.

HPC Innovation Lab. February 2018

Overview

Not long ago, [PowerEdge R740 Server](#) was released as part of Dell's 14th Generation server portfolio. It is a 2U Intel SkyLake based rack mount server and provides the ideal balance between storage, I/O and application acceleration. Besides VDI and Cloud, the server is also designed for HPC workloads. Compared to the previous R730 server, one of the major changes on GPU support is that, R740 supports up to 3 double width cards, which is one more than what R730 could support. This blog will focus on the performance of a single R740 server with 1, 2 and 3 Nvidia Tesla V100-PCIe GPUs. Multiple cards scaling number for HPC applications like [High-Performance Linpack \(HPL\)](#), [High Performance Conjugate Gradients benchmark \(HPCG\)](#) and [Large-scale Atomic/Molecular Massively Parallel Simulator \(LAMMPS\)](#) will be presented.

Table1: Details of R740 configuration and software version

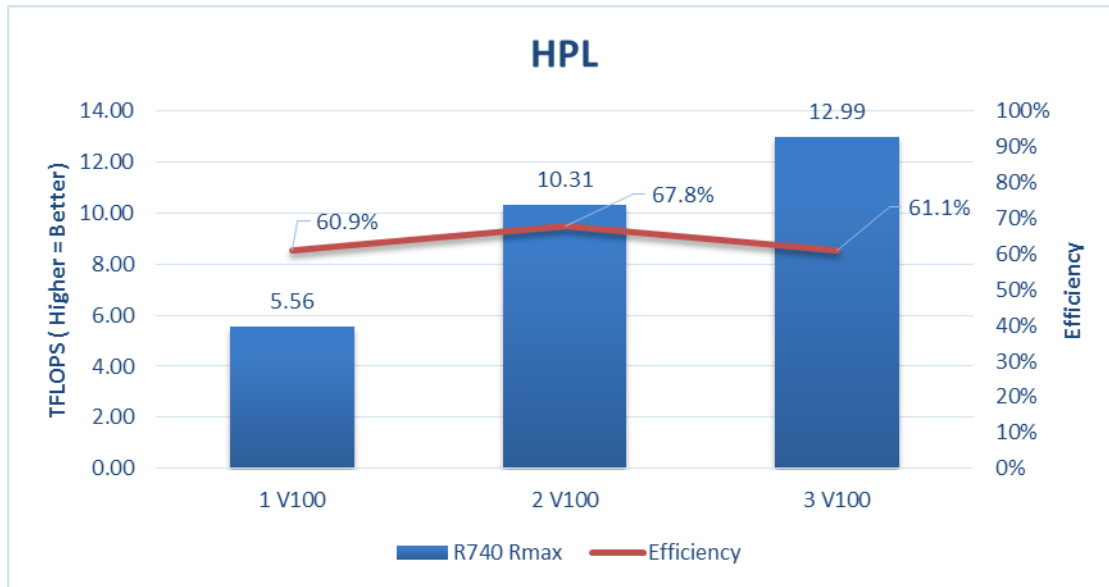
Processor	2 x Intel(R) Xeon(R) Gold 6150 @2.7GHz, 18c
Memory	384G(12*32G@2666MHz)
Local Disk	480G SSD
Operating System	Red Hat Enterprise Linux Server release 7.3
GPU	Nvidia Tesla V100-PCIe
CUDA Driver	387.26
CUDA Toolkit	9.1.85
Processor Settings > Logical Processors	Disabled
System Profiles	Performance

High Performance Linpack (HPL)

Figure 1: HPL Performance and efficiency with R740

Figure 1 shows HPL performance and efficiency numbers. Performance data increases with multiple GPUs nearly linearly. Efficiency line isn't flat, the peak 67.8% appears with the 2-card, which means configuration with 2 GPUs is the most optimized one for HPL. Number of 1 and 3 cards are about 7% lower than 2 card and they are affected by different factors:

- For the 1 card case, this GPU based HPL application designs to bond CPU and GPU. While running in large scale



with multiple GPU cards and nodes, it is known to make data access more efficient to bond GPU to the CPU. But testing with only 1 V100 here is a special case that only the 1st CPU bonded with the only GPU, and no workload on the 2nd CPU. Comparing with 2 cards result, the non-using part in [Rpeak](#) increased so efficiency dropped. This doesn't mean GPU less efficient, just because HPL application designs and optimizes for large scales, and 1 GPU only situation is special for HPL.

- For the 3 cards case, one of the major limitation factor is 3x1(P and Q) matrix need to be used. HPL is known to perform better with squared PxQ matrix. We also verified 2x2 and 4x1 Matrix on C4140 with 4x GPUs, and as a result 2x2 did perform better. But keep in mind, with 3 cards, the [Rmax](#) increased significantly as well. HPL is an extreme benchmark, in real world the capability of having the additional 3rd GPU will give a big advantage for efficiency with different application and different dataset.

HPCG

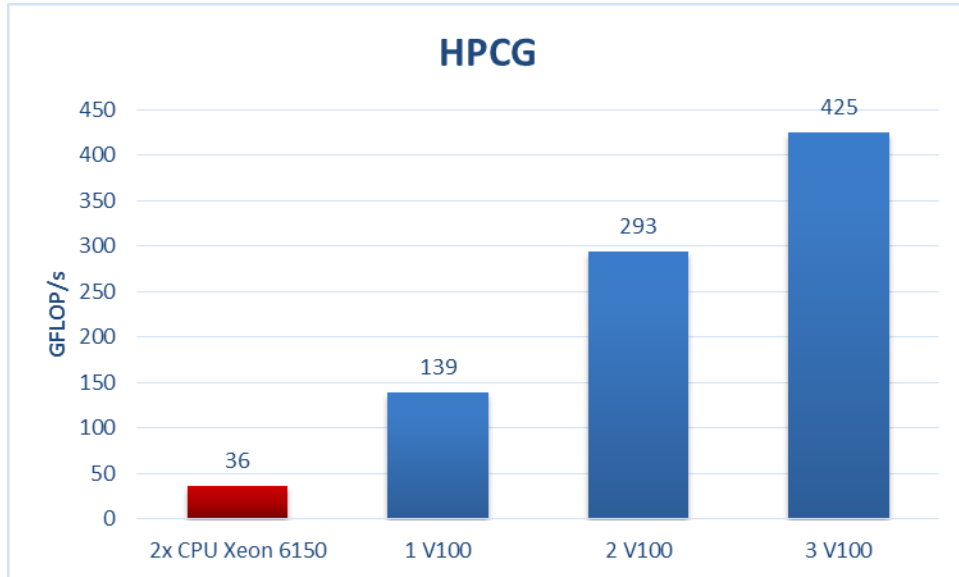


Figure 2: HPCG Performance with R740

As shown in Figure 2, comparing with dual Xeon 6150 CPU only performance, single V100 is already 3 times faster, 2 V100 is 8 times and 3 V100 (CPU only node vs GPU node) is nearly 12 times. Just like HPL, HPCG is also designed for a large scale, so single card performance isn't as efficient as multiple cards on HPCG. But unlike HPL, 3 card performance on HPCG is linearly scaled. It is 1.48 times higher than 2 cards', which is very close to the theoretical 1.5 times. This is because all the HPCG workload is run on GPU and its data fits in GPU memory. This proves application like HPCG can take the advantage of having the 3rd GPU.

LAMMPS

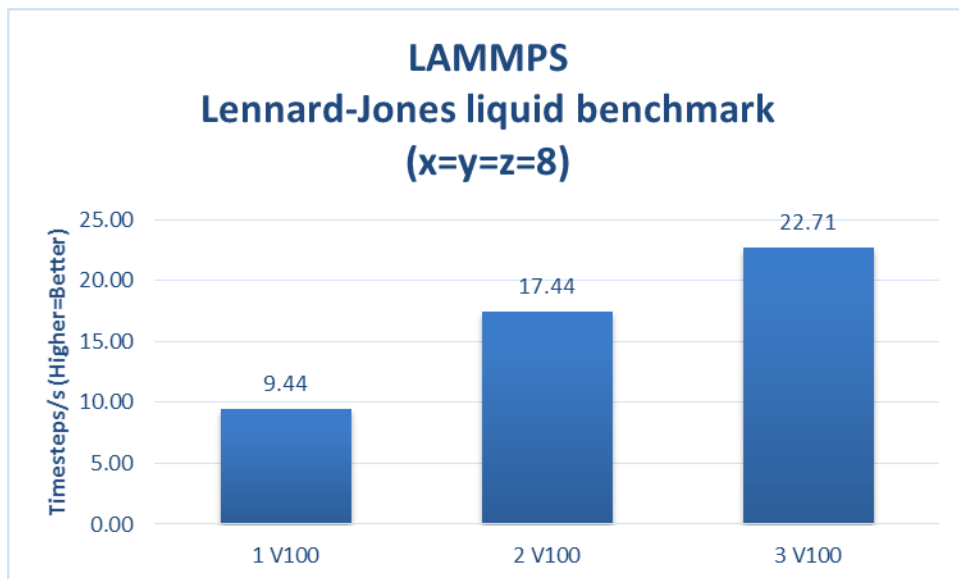


Figure 3: LAMMPS Performance with R740

The [LAMMPS](#) version used for this testing is 17Aug2017, which is the latest stable version at the time of testing. The testing dataset is in.intel.lj, which is the same one in all pervious GPU LAMMPS testing and it can be found [here](#). With the same parameters set from [previous testing](#), the initial values of space were $x=4$, $y=z=2$, the simulation executes with 512k atoms. Weak scaling obvious as timesteps/s number of 2 and 3 cards only 1.5 and 1.7 times than single card's. The reason for this is that the workload isn't heavy enough for 3 V100 GPUs. After adjusting all x,y,z to 8, 16M atoms generated in simulation, and then the performance scaled well with multiple cards. As shown in Figure 3, 2 and 3 cards is 1.8 and 2.4 times faster than single card, respectively. This results of LAMMPS is another example for GPU accelerated HPC applications that can benefit from having more GPUs in the system.

Conclusion

The R740 server with multiple Nvidia Tesla V100-PCIe GPUs demonstrates exceptional performance for applications like HPL, HPCG and LAMMPS. Besides balanced I/O, R740 has the flexibility for running HPC applications with 1, 2 or 3 GPUs. The newly added support for an additional 3rd GPU provides more compute power as well as larger total memory in GPU. Many applications work best when data fits in GPU memory and having the 3rd GPU allows fitting larger problems with R740.

References:

PowerEdge R740 Technical Guide: http://i.dell.com/sites/doccontent/shared-content/data-sheets/en/Documents/PowerEdge_R740_R740xd_Technical_Guide.pdf

