# Enhanced Molecular Dynamics Performance with K80 GPUs

By: Saeed Iqbal & Nishanth Dandapanthula

The advent of hardware accelerators in general has impacted Molecular Dynamics by reducing the time to results and therefore providing a tremendous boost in simulation capacity (E.g., previous NAMD blogs).  Over the course of time, applications from several domains including Molecular Dynamics have been optimized for GPUS.  A comprehensive (although a constantly growing) list can be found here. LAMMPS and GROMACS are two open source Molecular Dynamics (MD) applications which can take advantage of these hardware accelerators.

LAMMPS stands for "Large-scale Atomic/Molecular Massively Parallel Simulator" and can be used to model solid state materials and soft matter . GROMACS is short for "GROningen MAchine for Chemical Simulations". The primary usage for GROMACS is simulations for biochemical molecules (bonded interactions) but because of its efficiency in calculating non-bonded interactions (atoms not linked by covalent bonds), the user base is expanding to non-biological systems.

NVIDIA's K80 offers significant improvements over the previous model the K40.  From the HPC prospective *the most important improvement is the 1.87 TFLOPs (double precision) compute capacity*, which is about 30% more than K40.  The auto-boost feature in K80 automatically provides additional performance if additional power head room is available. The internal GPUs are based on the GK210 architecture and have a total of 4,992 cores which represent a 73% improvement over K40.  The K80 has a total memory of 24GBs which is divided equally between the two internal GPUs; this is a 100% more memory capacity compared to the K40.   The memory bandwidth in K80 is improved to 480 GB/s.  The rated power consumption of a single K80 card is a maximum of 300 watts.

Dell has introduced a new high density GPU server, PowerEdge C4130, it offers five configurations, noted here as "A" through "E".  Part of the goal of this blog is to find out which configuration is best suited for LAMMPS and GROMACS. The three quad GPU configurations "A", "B" and "C" are compared. Also the two dual GPU configurations "D" and "E" are compared for users interested in lower GPU density of 2 GPU per 1 rack unit.  The first two quad GPU configurations ("A" & "B") have an internal PCIe switch module which allows seamless peer to peer GPU communication. We also want to understand the impact of the switch module on LAMMPS and GROMACS. Figure 1 below shows the block diagrams for configurations A to E.

Combining K80s with the PowerEdge C4130, results in an extra-ordinarily powerful compute node. The C4130 can be configured with up to four K40 or K80 GPUs in a 1U form factor. Also the uniqueness  of PowerEdge C4130 is that it offers several workload specific configurations, potentially making it a better fit, for MD codes in general , and specifically for LAMMPS and GROMACS.
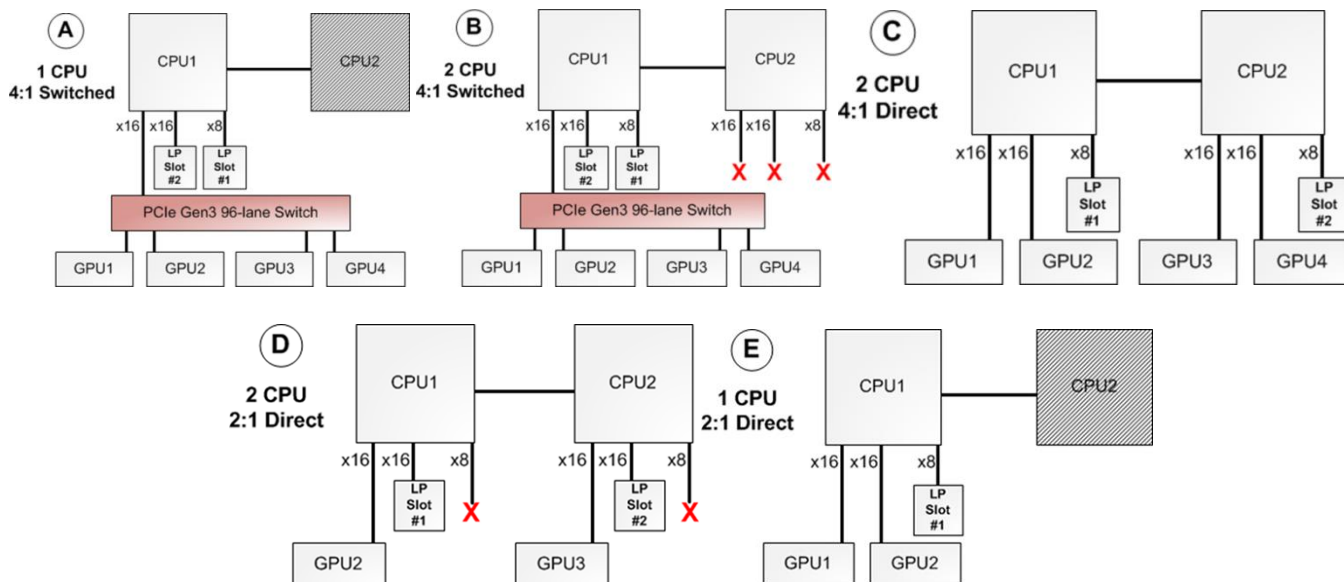
Figure 1: C4130 Configuration Block Diagram

Recently we have evaluated the performance of NVIDIA's Tesla K80 GPUs on Dell's PowerEdge C4130 server on standard benchmarks and applications (HPL and NAMD).

## Performance Evaluation with LAMMPS and GROMACS

In this blog, we quantify the performance of two of the molecular dynamics applications; LAMMPS and GROMACS by comparing their performance on K80s to a CPU only.   The performance is measured as  "Jobs/day" and "ns/day" (inverse of the number of days required to simulate 1 nanosecond of real-time) for LAMMPS and GROMACS respectively. Higher is better for both cases. Table 1 gives more information about the hardware configuration and application details used for the tests.

Table 1: Hardware Configuration and Application Details

| Server | PowerEdge C4130 |
|---|---|
| Processor | 1 or 2 x Intel Xeon CPU E5-2690 v3 @ 2.6 GHz (12 core) |
| Memory | 64GB or 128GB @ 2,133MHz |
| GPU | 2 or 4 x NVIDIA Tesla K80 (4,992 CUDA cores, base clock 562 MHz, boost clock 875MHz, power 300W) |
| Power supply | 2 x 1,600W |
| Operating System | RHEL 6.5 – kernel 2.6.32-431.el6.x86_64 |
| BIOS options | System Profile – Performance |
| | Logical Processor – Disabled |
| | Power Supply Redundancy Policy – Not Redundant |
| | Power supply Hot Spare – Disabled |
| | P2P – Enabled |

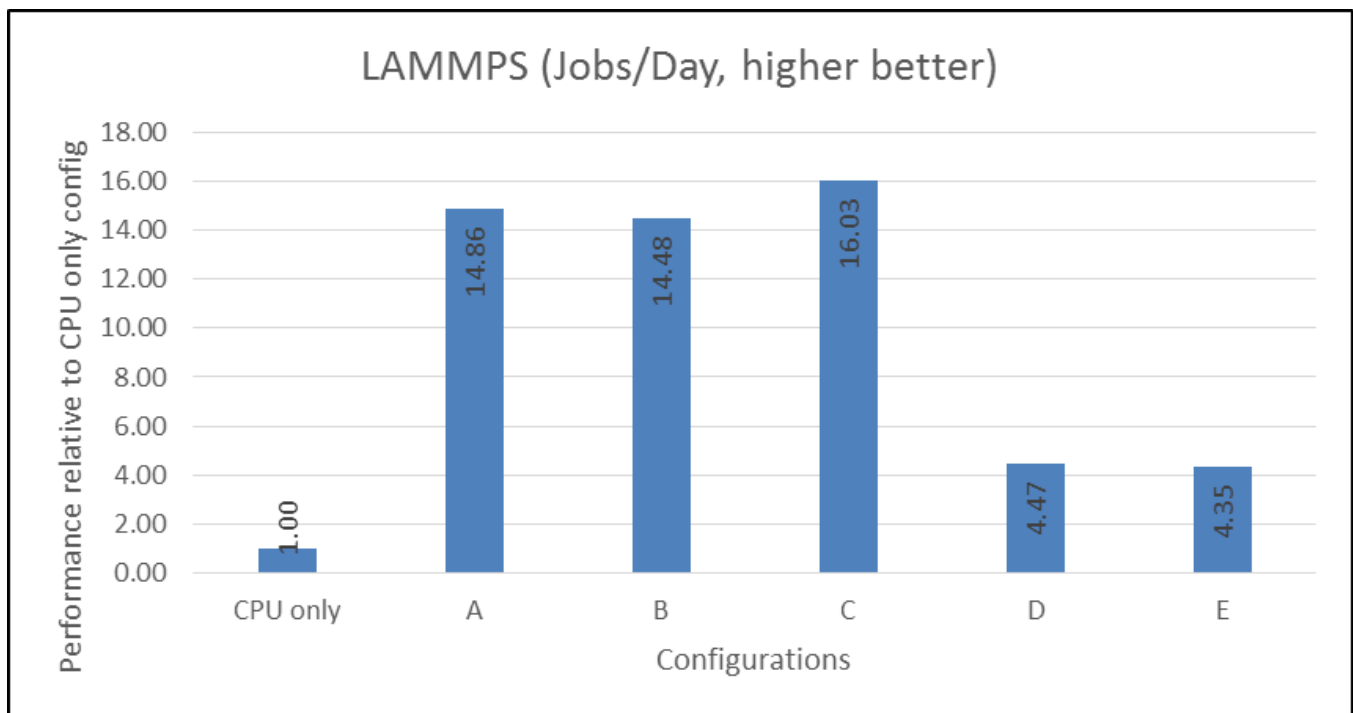| | Node Interleaving – Disabled |
|---|---|
| **CUDA Version and driver** | CUDA 6.5 (340.46) |
| **BIOS firmware** | 1.1.0 |
| **iDRAC firmware** | 2.02.01.01 |
| **LAMMPS** | 1 Feb 2014 stable version using lib/CUDA for GPU acceleration |
| | Benchmark: LJ (128 x 128 x 128) |
| **GROMACS** | 4.6.6 |
| | Benchmark: Water 0768 |



Figure 2: LAMMPS performance on K80s relative to CPUs

Figure 2 quantifies the performance of LAMMPS over the five configurations mentioned above and compares them to the CPU only server (CPU only server => performance of application on a server with two CPUS). The graph can be described as follows.

- Configurations A and B are the switched configurations with the only difference being that B has an extra CPU. Since LAMMPS just uses the GPU cores, the extra CPU does not offset the scale in terms of performance.

- Configurations "A", "B" and "C" are four GPU configurations.  Configuration C performs better than A and B. This can be attributed to the PCIe switch in configurations A and B which introduces an extra hop latency when compared to "C" which is a more balanced configuration.
- Among the two GPU configurations are D and E. Configuration D performs slightly better than E and this could again be attributed to the balanced nature of D. As mentioned previously, LAMMPS is not offset by the extra CPU in D.
- An interesting observation here is that when moving from 2 K80s to 4 K80s (i.e. comparing D and C configurations in  Figure 2) the performance almost quadruples. This shows that for each extra K80 added (2 GPUs per K80) the performance doubles. This can be partially attributed to the size of the dataset used.
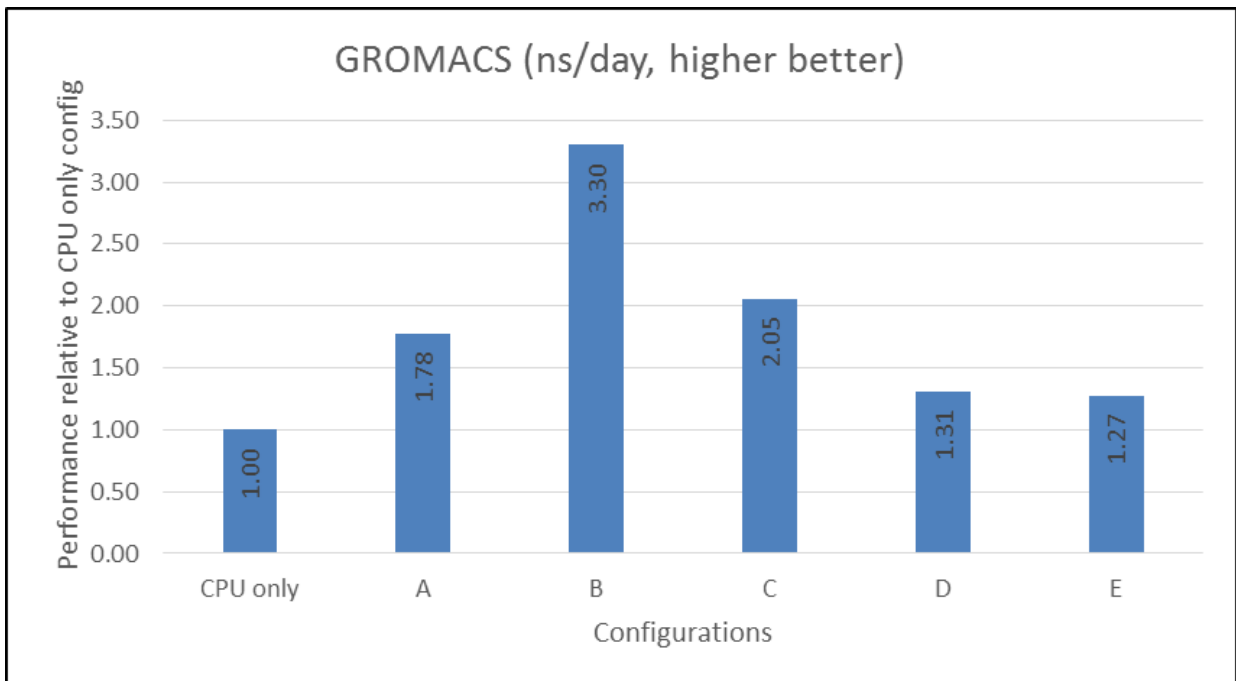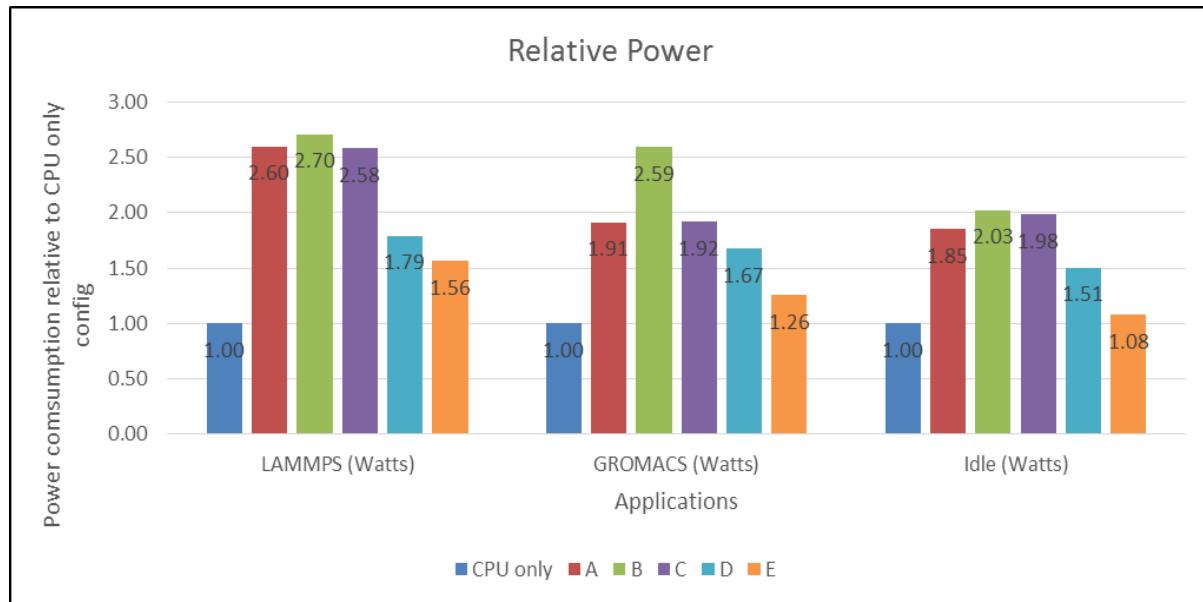


Figure 3: GROMACS performance on K80s relative to CPUs

Figure 3 shows the performance of GROMACS among the five configurations and the CPU only configuration. The explanation is as follows.

- Among the quad CPU configurations (A, B and C), B performs the best. In addition to the 4 GPUs attached to CPU1, GROMACS also used the whole second CPU2 making B the best performing configuration.    It seems GROMACS benefits from the second CPU as well as the switch, it's likely that application has substantial GPU to GPU communication.
- Configuration C outperforms A. This can be attributed to the more balanced nature of C. Another contributing factor may be the latency hit because of the PCIe switch in A.
- Even in the dual GPU configurations (D and E), D which is the balanced of the both, slightly outperforms E.

Performance is not the only criteria when a performance optimized server as dense as the Dell PowerEdge C4130 with 4 x 300 Watt accelerators is used. The other dominating factor is how much power these platforms consume. Figures 4 answers questions pertaining to power.

- In case of LAMMPS the order of power consumption is as follows. B > A >= C > D > E
    - Configuration B is a switched configuration and has an extra CPU then Configuration A.
    - Configuration A incurs a slight overhead of the switch and thus takes up slightly more power than C.
    - Configuration D is a dual GPU, dual CPU configuration and thus takes up more power than E, which is a single CPU dual GPU configuration

- In case of GROMACS, the order is still the same, but the B takes up considerably more power than A and C when compared to LAMMPS. This is because GROMACS uses the extra CPU in B while LAMMPS does not.



In conclusion, Both GROMACS and LAMMPS benefit greatly from Dell's PowerEdge C4130 servers and NVIDIAs K80s. *In the case of LAMMPS, we see a 16x improvement with only a 2.6x more power. In case of GROMCAS, we see a 3.3x improvement in performance while talking up 2.6x more power.* The comparisons in this case are with a dual CPU only configuration. Obviously, there are a lot other factors which come into play when scaling these results to multiple nodes; GPU direct, interconnect, size of the dataset/simulation are just a few of those.