DRAGNET Cluster Benchmark Numbers

Parts of the cluster and interconnects have been stress tested to optimize configuration and to find upper performance bounds that can be useful for application optimization and rough capacity estimates.

The tests described tend to cover *maximum* achievable performance results on a synthetic ideal workload. It is very likely that your (real) application will never reach these numbers, as the workload is non-ideal, and reaching peak performance can take a lot of effort.

Cluster Specifications

Networking

We have benchmarked the infiniband and 10G networks. A good guide is available at the http://fasterdata.es.net/ under Host Tuning (and to a lesser extend under Network Tuning). But the indicated Linux kernel sysctl knobs did not help; CentOS 7 already has decent settings, and our transfers are all on a low latency LAN (as opposed to wide-area).

Infiniband

Each drgXX node has an FDR (54.545 Gbit/s) HCA (Host Channel Adapter). The 36 port cluster switch is connected to the COBALT switch with 5 aggregated lines (272.727 Gbit/s). See below what can be achieved under ideal circumstances.

IPolB: TCP and UDP

An application that uses the Infiniband (ib) network normally uses IPoIB (IP-over-Infiniband) to transfer data via TCP or UDP. DRAGNET IPoIB settings have been optimized for TCP (at the cost of UDP performance). We (mostly) use TCP and will not receive UDP data from LOFAR stations directly.

We used the iperf3 benchmark and got the following bandwidth numbers between two drgXX nodes:

```
Out-of-the-box TCP bandwidth: 26-28 Gbit/s. We can ''get iperf3'' TCP bw to 45.4 Gbit/s:

# Set CPU scaling gov to 'performance' (default is 'powersave') (from 39.3 to 45.4 Gbit/s (TCP, 2 streams))

$ for i in /sys/devices/system/cpu/cpu*/cpufreq/scaling_governor; do echo performance | sudo tee $i; done

[amesfoort@drg23 ~]$ sudo iperf3 -A 10 -B drg23-ib -s
[amesfoort@drg21 ~]$ sudo iperf3 -N -4 -A 10 -B drg21-ib -i 1 -l 1M -P 2 -t
```

```
20 -c drg23-ib

-A 10: sets CPU affinity to hw thread 10. mlx_4 is on the 2nd CPU (8-15(,24-31)) and the INT handler happens to be on hw thr 9
-A 10 (or another suitable nr): from 26.2 to 45.4 Gbit/s (tcp, 2 streams)

-P 2: 2 parallel streams instead of 1: from 36.4 to 45.4 Gbit/s (tcp, -A 10)
```

I have extensively tried the sysctl knobs of the Linux kernel networking settings (net.core.*, net.ipv4.*, txqueuelen), but I did not get an improvement. Likely, Linux kernel autotuning + CentOS 7 settings are already ok for local area networking up to at least 45 Gbit/s.

A real application (not a synthetic benchmark) likely does something else except for data transfer and may have trouble reaching these numbers, because CPU load is a limiting factor and the clock frequency boost of 1 core on DRAGNET CPUs is higher (up to 3.2 GHz) than for all cores at once (up to 2.6 GHz). Standard DRAGNET CPU clock frequency is 2.4 GHz.

Numbers (before tuning iperf3 tests) obtained using qperf can be seen below under the RDMA sub-section.

RDMA

RDMA (Remote Direct Memory Access) allows an application to directly access memory on another node. Although some initial administration is set up via the OS kernel, the actual transfer commands and completion handling does not go via the kernel. This also saves data copying on sender and receiver and CPU usage.

Typical applications that may use RDMA are applications that use MPI (Message Passing Interface) (such as COBALT), or (hopefully) the LUSTRE client. NFS can also be set up to use RDMA. You can program directly into the verbs and rdma-cm C APIs and link to those libraries, but be aware that extending some code to do this is not a 1 hr task... (Undoubtedly, there is also a Python module that either only wraps or even makes makes life easier.)

We used the qperf benchmark and got the following bandwidth and latency numbers between two drgXX nodes (TCP/UDP/SCTP over IP also included, but not as fast as mentioned above):

```
[amesfoort@drg22 ~]$ qperf drg23-ib sctp_bw sctp_lat tcp_bw tcp_lat udp_bw
udp_lat
sctp_bw:
    bw = 355 MB/sec
sctp_lat:
    latency = 8.94 us
tcp_bw:
    bw = 3.65 GB/sec
tcp_lat:
    latency = 6.33 us
udp_bw:
    send_bw = 6.12 GB/sec
    recv_bw = 3.71 GB/sec
udp_lat:
```

```
latency = 6.26 us
```

```
[amesfoort@drg22 ~]$ sudo qperf drg23-ib rc bi bw rc bw rc lat uc bi bw
uc bw uc lat ud bi bw ud bw ud lat
rc bi bw:
    bw = 11.9 GB/sec
rc bw:
    bw = 6.38 GB/sec
rc lat:
    latency = 5.73 us
uc bi bw:
    send bw = 12 GB/sec
    recv bw = 11.9 GB/sec
uc bw:
    send_bw = 6.24 GB/sec
    recv bw = 6.21 \, \text{GB/sec}
uc lat:
   latency = 4.03 us
ud bi bw:
    send bw = 10.1 \text{ GB/sec}
    recv bw = 10.1 GB/sec
ud bw:
    send bw = 5.93 \text{ GB/sec}
    recv bw = 5.93 \text{ GB/sec}
ud lat:
    latency = 3.94 us
```

```
[amesfoort@drg22 ~]$ sudo qperf drg23-ib rc rdma read bw rc rdma read lat
rc_rdma_write_bw rc_rdma_write_lat rc_rdma_write_poll_lat uc_rdma_write_bw
uc rdma write lat uc rdma write poll lat
rc rdma read bw:
   bw = 5.6 GB/sec
rc_rdma_read_lat:
   latency = 4.99 us
rc rdma write bw:
   bw = 6.38 GB/sec
rc rdma write lat:
   latency = 5.35 us
rc rdma write poll lat:
   latency = 925 ns
uc rdma write bw:
   send bw = 6.26 \text{ GB/sec}
    recv bw = 6.23 GB/sec
uc_rdma_write_lat:
   latency = 3.58 us
uc_rdma_write_poll_lat:
   latency = 922 ns
```

[amesfoort@drg22 ~]\$ sudo qperf drg23-ib rc_compare_swap_mr rc_fetch_add_mr
ver_rc_compare_swap ver_rc_fetch_add
rc_compare_swap_mr:

```
msg rate = 2.08 M/sec
rc_fetch_add_mr:
   msg_rate = 2.14 \text{ M/sec}
ver_rc_compare_swap:
   msg_rate = 2.1 M/sec
ver_rc_fetch_add:
   msg rate = 2.4 M/sec
```

From:

https://www.astron.nl/lofarwiki/ - LOFAR Wiki

Permanent link:

https://www.astron.nl/lofarwiki/doku.php?id=dragnet:cluster_benchmark&rev=1469642248



