

Exploring the internal messaging performance of MQ as part of CP4I deployed on OCP 4.2

Objective

To illustrate the performance capability of MQ when deployed on OpenShift Cloud Platform (OCP) 4.2 as part of Cloud Pak for Integration (CP4I).

Environment

A bare metal setup comprising 3 master nodes, 1 bastion node and 9 worker nodes, although for this test only 2 worker nodes are used (1 for the QM and 1 for the MQ Client).

OpenShift Cloud Platform Version: 4.2.29

CP4I Version: 2019.4.1

The persistent storage uses a Fibre channel volume to a remote SAN.

The MQ version for that level of CP4I is usually MQ 9.1.3 (https://www.ibm.com/support/knowledgecenter/en/SSFKSJ_9.1.0/com.ibm.mq.ctr.doc/ctr_supported_versions.htm). The image for these tests contains a modified image containing MQ 9.1.4 and a set of performance improvements (see next section).

The default cluster SDN (Software Defined Network) utilizes 1GbE networking. All master and worker nodes are also connected by an additional 10GbE network. The Multus additional network support has been used to allow the client and QM to communicate over the 10GbE network, please see separate guidance (<https://github.com/ibm-messaging/mqperf/blob/gh-pages/openshift/configuration.md>) on how this was setup.

The number of threads supported in a container in this environment is currently limited to 1024. To be able to run with more threads, a configuration change to the cri-o.conf file is required. cri-o is an implementation of the Container Runtime Interface (CRI) that supports Open Container Initiative (OCI) compatible runtimes. Again, please see the separate guidance on OpenShift configuration on how this was setup.

The first two sections in this report use a CPU limit of 32 cores for the QM. The CPU limit for the client is set to a value to avoid the client encountering CPU starvation.

Please see Appendix A for the specification of the hardware used for the cluster nodes.

Changes to MQ image

Three changes were made to the base MQ image to improve performance. These will likely be included (or at least configurable) in future CP4I releases.

- Enable FASTPATH bindings
- Increase number of MaxChannels/MaxActiveChannels to 999999999
- Increased log configuration to 64 primary files (of 16384 4K Pages) resulting in 4GB log allocation

Scenario

The scenario that will be used in the testing for this whitepaper is the standard requester/responder scenario as featured in our distributed performance reports.

The MQ client runs in its own container with a fixed number of responders (500) connecting to the QM under test. The test then iterates through an increasing number of client requesters which sends messages across 10 request queues. The responders consume the messages from the request queues and place them on the reply queues where the requester clients obtain their specific reply (via correlation ID) to their original message.

A full round trip is 2 message puts and 2 message gets. The client runs within the OCP cluster and connects to the QM using the address allocated to the QM Pod on the additional 10GbE network and port 1414.

For this investigation, 2KB, 20KB and 200KB persistent messages are used. TLS is not used as all messaging data flows over the private 10GbE network and the MQ client and QM are both running in the same OpenShift cluster.

Non Persistent Results

The graph below shows how the MQ QM performs for a 2K message size.

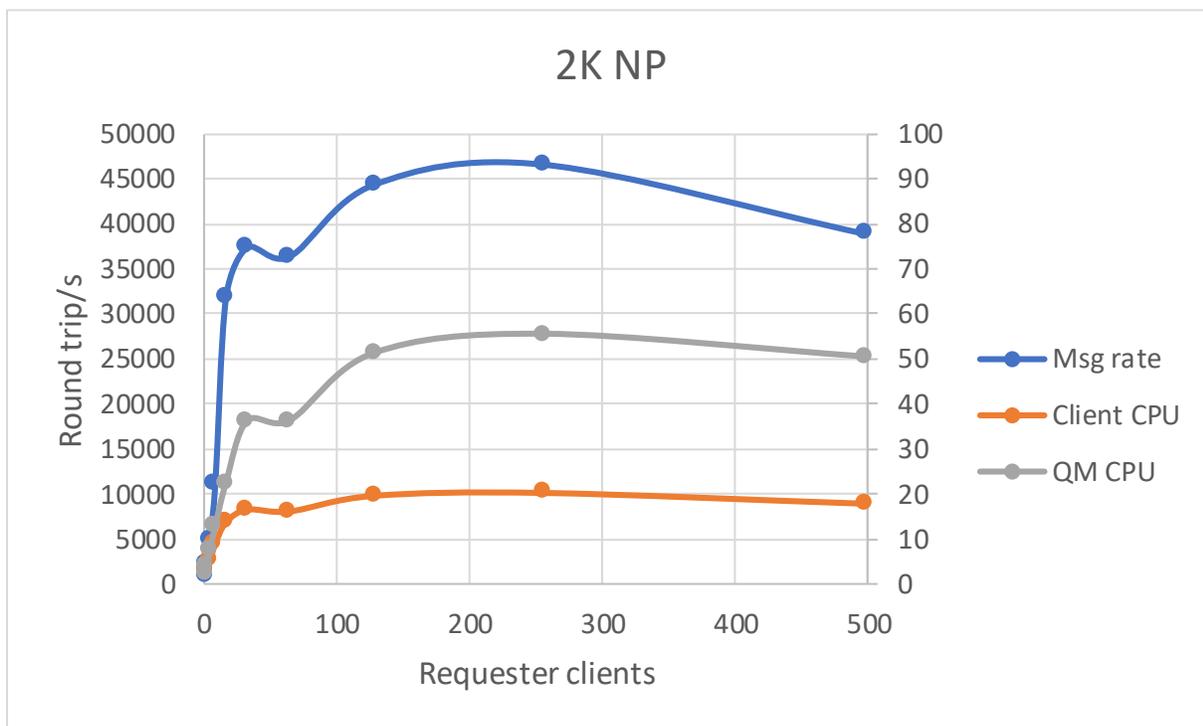


Figure 1 – 2K Non Persistent

Note that the reported CPU is based on the full capacity of the worker node, which in this case is 64 Hyperthreaded cores. So a pod restricted to 32 cores would report as having used up to 50% of the available capacity. There are additional pods running on that Node to support the management and configuration of the OCP cluster which is why the maximum reported value is approximately 55%

The above graph shows that the QM can achieve a peak throughput of over 45,000 round trips/s, and even with just 16 requester client threads, the QM can achieve over 30,000 round trips/s.

The graph below shows how the MQ QM performs for a 20K message size.

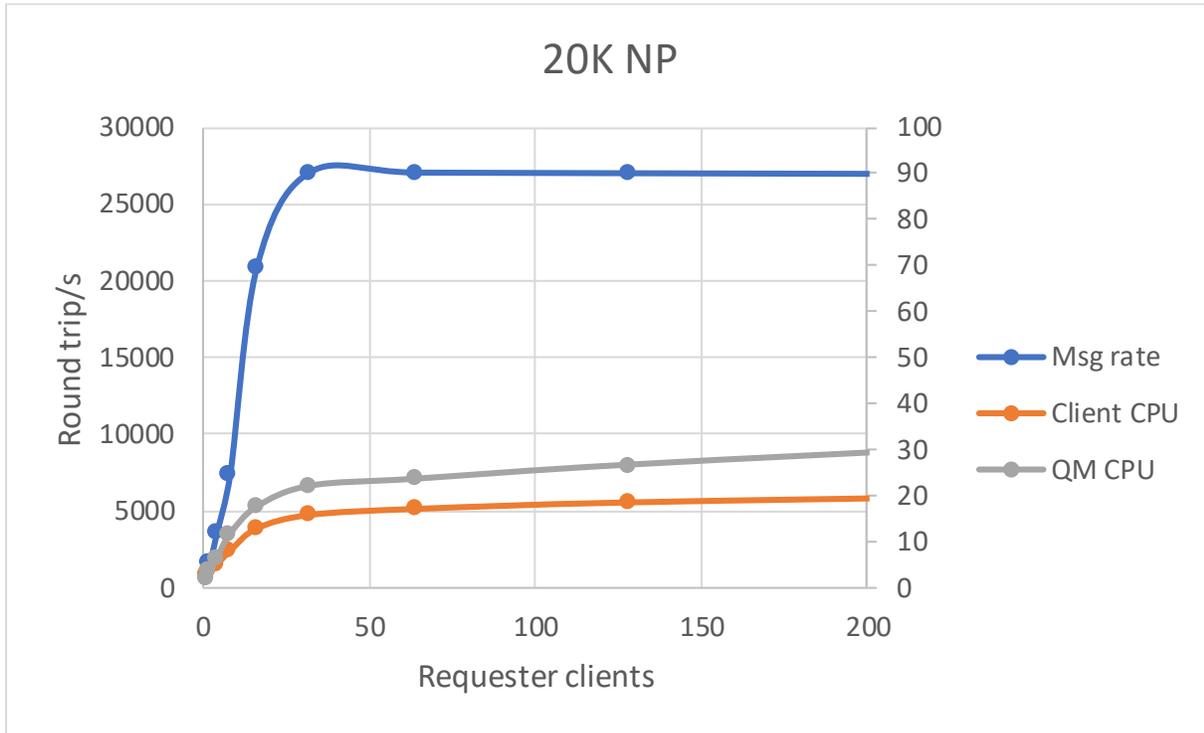


Figure 2 - 20K Non Persistent

The above graph shows that as we increase the message size to 20K, the QM CPU is no longer the limiting factor and we are now limited by our 10GbE network at over 27,000 round trips/sec, achievable with 32 or more requester clients.

The graph below shows how the MQ QM performs for a 200K message size.

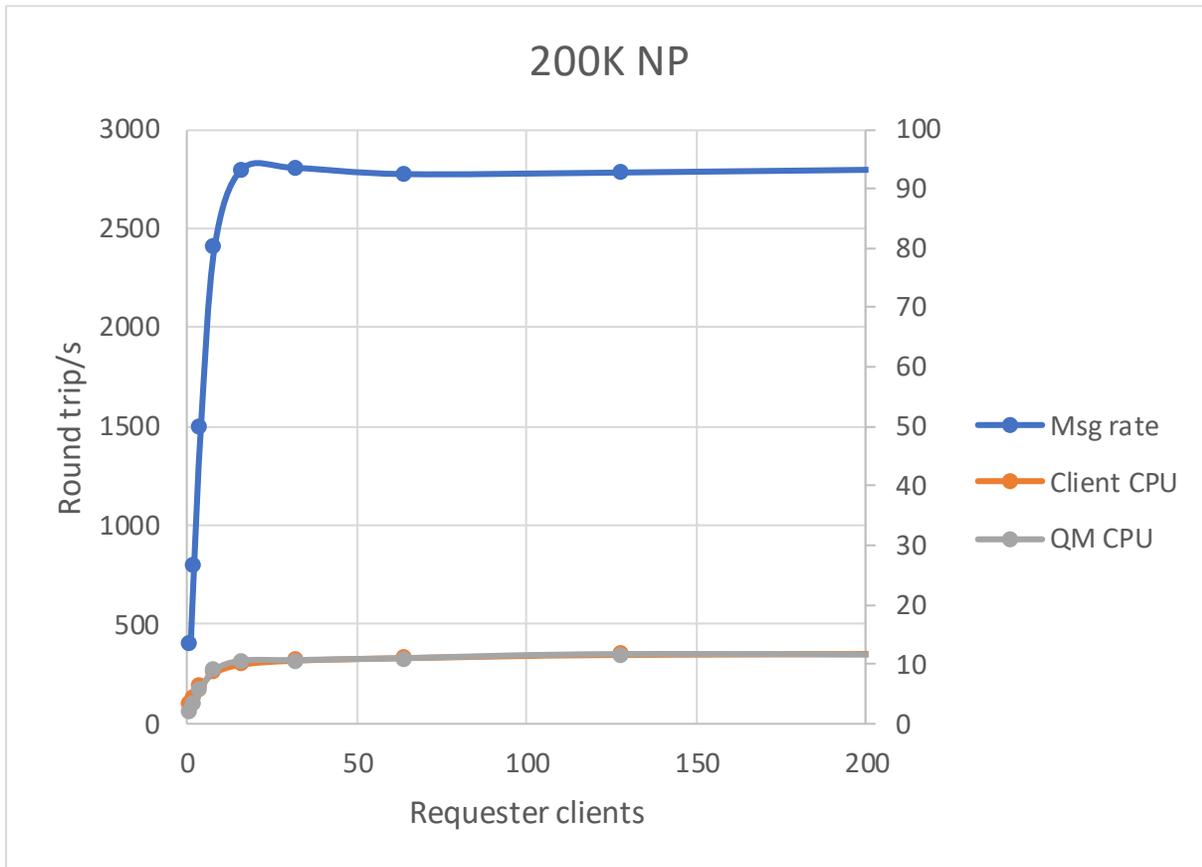


Figure 3 - 200K Non Persistent

The above graph again shows that we are limited by the network when the throughput has reached over 2,700 round trips/s from 16 threads and the QM CPU is barely over 10% utilised.

Persistent Results

The graph below shows how the MQ QM performs for a 2K message size.

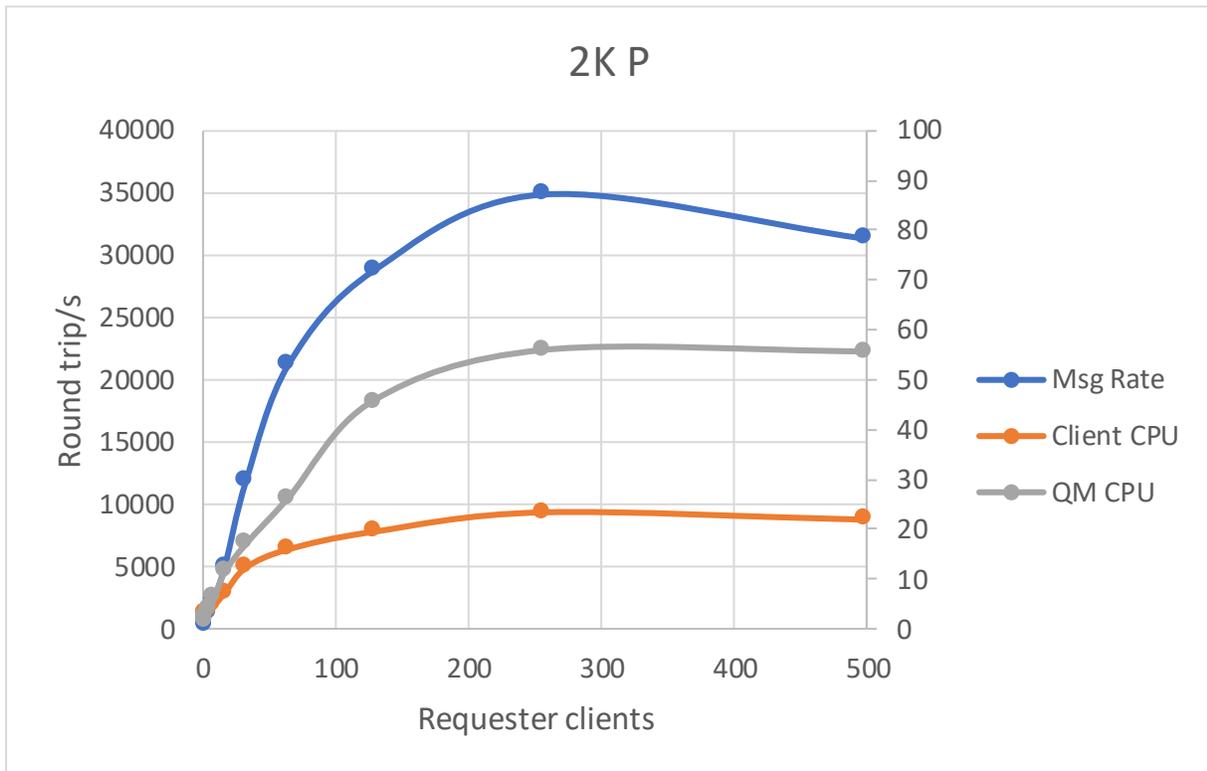


Figure 4 - 2K Persistent

The above graph shows that the QM can achieve a peak throughput of nearly 35,000 round trips/s until we saturate all of the CPU available to the QM which has a CPU limit of 32 cores.

The graph below shows how the MQ QM performs for a 20K message size.



Figure 5 - 20K Persistent

The above graph shows that as we increase the message size to 20K, the QM CPU is no longer the limiting factor and we are now limited by the persistence layer at over 8,000 round trips/sec

The graph below shows how the MQ QM performs for a 200K message size.

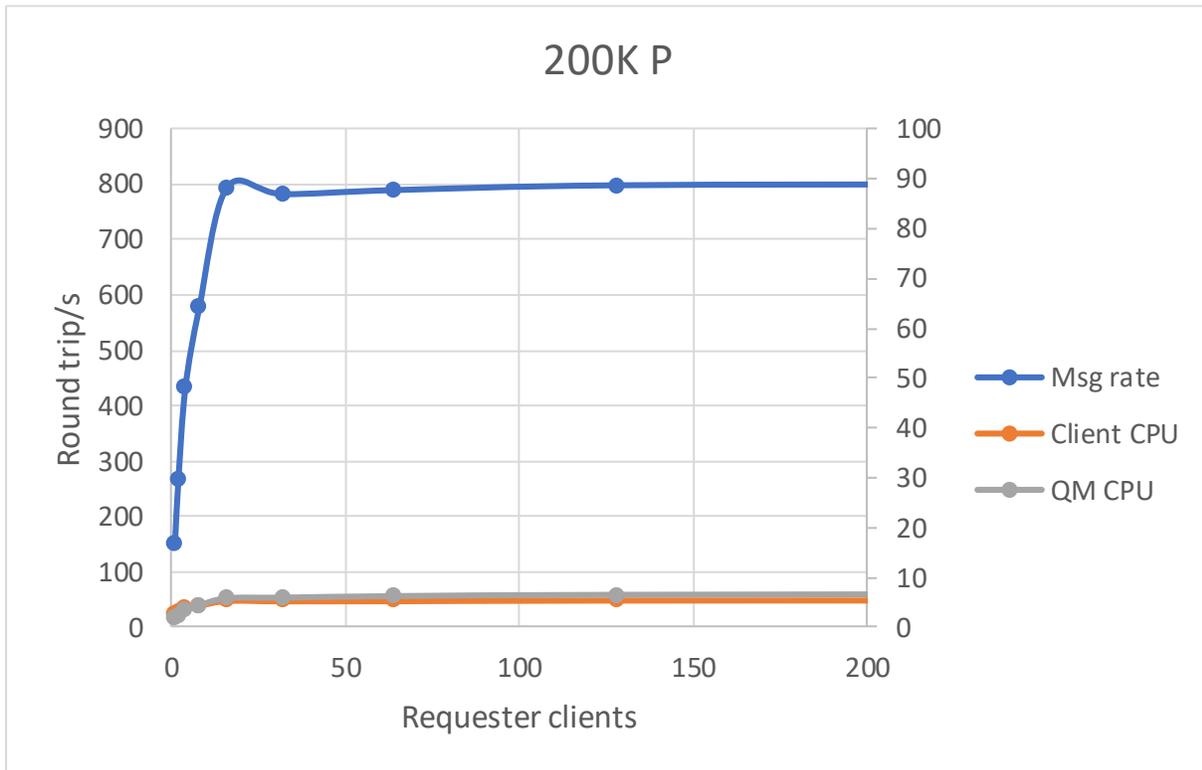


Figure 6 - 200K Persistent

The above graph again shows that we are again limited by the persistence layer when the throughput has reached over 800 round trips/s and the QM CPU is less than 10% utilised.

Scaling Results

The results presented so far have been with a CPU limit of 32 cores, which is greater than we expect the majority of scenarios to use. To illustrate how MQ performs in the OpenShift environment with varying levels of CPU resources, the 2K, 20K and 200K message tests have been run against the MQ QM in multiple CPU configurations and the peak throughput noted.

The graph below shows how the MQ QM scales across varying CPU cores with a 2K message size.

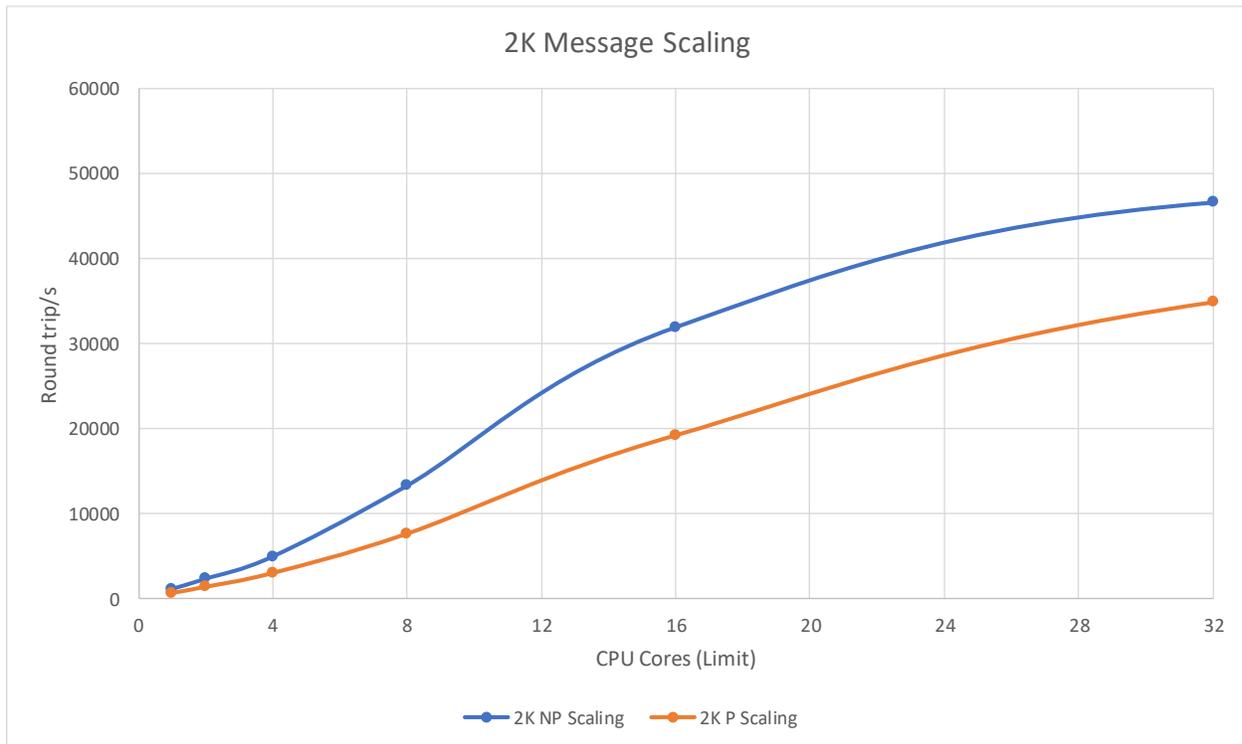


Figure 7 - 2K Scaling

The chart above shows how we can support over 1,000 round trips/s at a single CPU core right up to over 45,000 round trips/s at 32 CPU cores for Non Persistent messaging. For Persistent messaging the respective values are approximately 700 and 35,000 round trips/s.

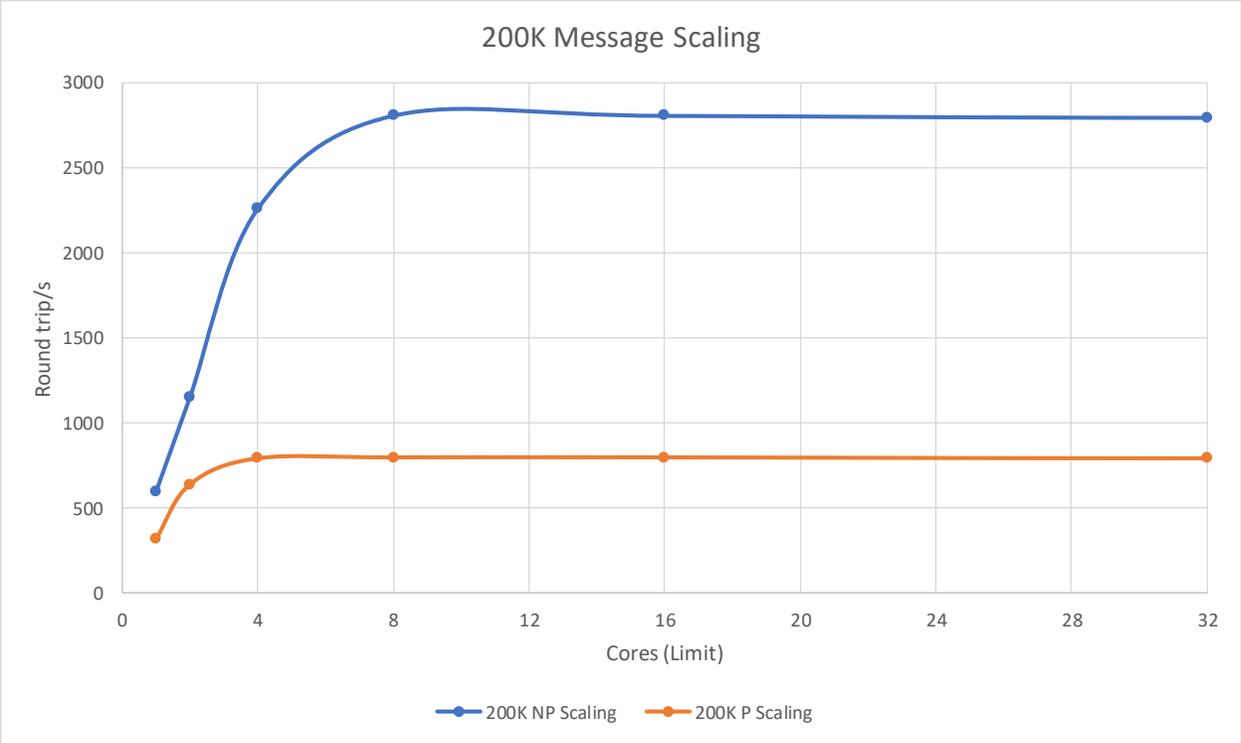
The graph below shows how the MQ QM scales across varying CPU cores with a 20K message size.



Figure 8 - 20K Scaling

The chart above shows how we can support 1,000 round trips/s at a single CPU core right up to over 25,000 round trips/s at 16 CPU cores for Non Persistent messaging. For Persistent messaging the respective values are approximately 600 and 8,000 round trips/s.

The graph below shows how the MQ QM scales across varying CPU cores with a 200K message size.



The chart above shows how we can support 600 round trips/s at a single CPU core right up to 2,800 round trips/s at 8 CPU cores for Non Persistent messaging. For Persistent messaging the respective values are approximately 300 and 800 round trips/s.

Conclusions

In this whitepaper we have looked at the performance of the MQ QM in the OpenShift environment and shown the affect of varying message size, requester clients and CPU cores have on the performance of the QM.

This data should help you size your solutions to support your intended workload. We will be producing further reports examining the impact of locating clients outside of the cluster and utilizing TLS to secure the messaging payload.

Appendix A

Hardware specification for Worker Nodes:

System	ThinkSystem SR630
CPU	2x16 Core 2.8Ghz Xeon Gold 6242 Hyperthreaded
RAM	96GB RAM RDIMM TruDDR4 2933MHz
RAID	930-16i 4GB Flash PCI 12Gb RAID Adapter
Disks	800GB SSD (2x400GB) SS530 Performance SAS 12Gbp/s
SAN Connectivity	Dual Port HBA 16Gb
10GbE Network	Dual Port 10GbE Broadcom Network Adapter
100GbE Network	Dual Port 100GbE Mellanox ConnectX-4 Network Adapter

<https://lenovopress.com/lp1049-thinksystem-sr630-server-xeon-sp-gen2>

Hardware specification for Master, Infrastructure and Bootstrap nodes:

System	ThinkSystem SR530
CPU	1x8 Core 2.1Ghz Xeon Silver 4208 Hyperthreaded
RAM	32GB RAM (2x16GB) RDIMM TruDDR4 2666MHz
RAID	530-8i PCI 12Gb RAID Adapter
Disks	480GB SSD (2x240GB) S4610 Mainstream SATA 6Gbp/s
10GbE Network	Dual Port 10GbE Broadcom Network Adapter

<https://lenovopress.com/lp1045-thinksystem-sr530-server-xeon-sp-gen2>