

# IBM MQ Appliance Performance Report

**Model: M2002**

**Version 2.0 - November 2020**

Sam Massey  
IBM MQ Performance  
IBM UK Laboratories  
Hursley Park  
Winchester  
Hampshire



# 1 Notices

## **Please take Note!**

Before using this report, please be sure to read the paragraphs on “disclaimers”, “warranty and liability exclusion”, “errors and omissions”, and the other general information paragraphs in the "Notices" section below.

## **First Edition, October 2018.**

This edition applies to *IBM MQ Appliance* (and to all subsequent releases and modifications until otherwise indicated in new editions).

© Copyright International Business Machines Corporation 2018, 2020. All rights reserved.

## Note to U.S. Government Users

Documentation related to restricted rights.

Use, duplication or disclosure is subject to restrictions set forth in GSA ADP Schedule contract with IBM Corp.

## **DISCLAIMERS**

The performance data contained in this report was measured in a controlled environment. Results obtained in other environments may vary significantly.

You should not assume that the information contained in this report has been submitted to any formal testing by IBM.

Any use of this information and implementation of any of the techniques are the responsibility of the licensed user. Much depends on the ability of the licensed user to evaluate the data and to project the results into their own operational environment.

## **WARRANTY AND LIABILITY EXCLUSION**

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law:

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION “AS IS” WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR A PARTICULAR PURPOSE.

Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore this statement may not apply to you.

In Germany and Austria, notwithstanding the above exclusions, IBM's warranty and liability are governed only by the respective terms applicable for Germany and Austria in the corresponding IBM program license agreement(s).

## **ERRORS AND OMISSIONS**

The information set forth in this report could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; any such change will be incorporated in new editions of the information. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this information at any time and without notice.

## **INTENDED AUDIENCE**

This report is intended for architects, systems programmers, analysts and programmers wanting to understand the performance characteristics of *IBM MQ Appliance*. The information is not intended as the specification of any programming interface that is provided by IBM. It is assumed that the reader is familiar with the concepts and operation of IBM MQ Appliance.

## **LOCAL AVAILABILITY**

References in this report to IBM products or programs do not imply that IBM intends to make these available in all countries in which IBM operates. Consult your local IBM representative for information on the products and services currently available in your area.

## **ALTERNATIVE PRODUCTS AND SERVICES**

Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

## **USE OF INFORMATION PROVIDED BY YOU**

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

## **TRADEMARKS AND SERVICE MARKS**

The following terms used in this publication are trademarks of their respective companies in the United States, other countries or both:

- **IBM Corporation** : IBM
- **Oracle Corporation** : Java

Other company, product, and service names may be trademarks or service marks of others.

## **EXPORT REGULATIONS**

You agree to comply with all applicable export and import laws and regulations.

## 2 Contents

1	Notices .....	2
2	Contents .....	4
3	Introduction .....	5
4	Request/Responder Scenario .....	7
4.1	Test Scenario C1 – 10 Applications per QM, 1 QM, Non-persistent .....	9
4.2	Test Scenario C2 – 10 applications per QM, 1 QM, Persistent.....	10
4.3	Test Scenario C3 – 10 applications per QM, 10 QM, Non-persistent.....	11
4.4	Test Scenario C4 – 10 applications per QM, 10 QM, Persistent.....	12
5	Connection Scaling .....	13
5.1	Connection Test .....	13
6	HA Scenarios .....	14
6.1	Test Scenario HA1 – 10 Applications per QM, 1 QM, Persistent.....	15
6.2	Test Scenario HA2 – 10 applications per QM, 10 QM, Persistent .....	16
6.3	How does HA perform over larger distances? .....	17
7	DR Scenarios .....	20
7.1	Test Scenario DR1 – 10 Applications per QM, 1 QM, Persistent.....	21
7.2	Test Scenario DR2 – 10 Applications per QM, 10 QM, Persistent.....	22
7.3	How does DR perform over larger distances?.....	23
8	HA and DR Scenarios .....	24
9	Additional M2002A vs M2002B scenarios .....	25
10	TLS .....	31
11	AMS .....	32
12	Frequently Asked Questions.....	34
13	Appendix A – Client machine specification .....	35
14	Appendix B – QM Configuration .....	35

### 3 Introduction

This performance report at version 2.0 contains performance data based on the MQ Appliance models M2002A and M2002B. This is an updated report featuring MQ 9.2. Whilst many scenarios will offer similar performance to MQ 9.1 as featured in version 1.0 of this report, there are some improvements that will be discussed. Please also see the first version of this report for a comparison with the preceding appliance M2001. This report covers standalone, TLS, AMS, HA and DR messaging performance and includes the following highlights:

- The performance of a single Non HA QM has been improved by up to 20%. See section 4.2
- The performance of a single HA QM has been improved by up to 80%. See section 6.1
- Over 150,000 round trips/second achieved in a TLS encrypted messaging scenario. See section 10
- Over 20,000 round trips/second achieved in an AMS encrypted messaging scenario. See section 11
- Nearly 110,000 round trips/second peak messaging rate achieved in an HA enabled scenario (~220,000 messages produced and ~220,000 messages consumed). See section 6.2
- Over 285,000 round trips/second peak messaging rate achieved in a NonPersistent messaging scenario (~570,000 messages produced and ~570,000 messages consumed). See section 4.1

The M2002 hardware components and how they compare to the previous model M2001 are shown below:

Model	M2001A	M2002A	M2002B
CPU	2x10 Core HT	2x12 Core HT	1x6 Core HT
RAM	192GB	192GB	192GB
Storage	3.2TB	6.4TB	3.2TB
IO Subsystem	RAID 1	RAID 10	RAID 10
Workload and replication network connectivity	4x10Gb 8x1Gb	6x10Gb 8x1Gb 4x40Gb	6x10Gb 8x1Gb 4x40Gb
Management	2x1Gb	2x1Gb	2x1Gb
Chipset	Ivybridge	Skylake	Skylake
RAID	6Gb/s 1GB cache	12Gb/s 2GB cache	12Gb/s 2GB cache

The MQ appliance combines all of the core MQ functionality with the convenience, ease of install and simplified maintenance of an appliance.

There are local disks within the appliance to enable efficient persistent messaging by the local Queue Managers. The four 3.2TB SSD drives are configured in a RAID10 configuration so that data is protected should one of the drives suffer a failure. High Availability (HA) may be achieved by the pairing of two MQ appliances which results in the Queue Manager (QM) log and queue files being distributed synchronously across the pair of appliances. Disaster Recovery (DR) may be achieved by the addition of a remote appliance to which QM data is distributed asynchronously. This report will also illustrate the HA and DR capabilities of the new model.

The MQ appliance can be purchased in two variants: M2002A and M2002B. There are two main differences for the M2002B as highlighted in the table above, reduced CPU capacity and reduced filesystem storage space.

As before, you can purchase an upgrade to convert an M2002B appliance to an M2002B+ appliance, which has the same capacity as an M2002A appliance.

The majority of the tests use the M2002A variant of the MQ Appliance and this is the default hardware unless stated otherwise. A number of tests were also conducted using the M2002B variant and provide comparative data points to the main testing to provide appropriate capacity planning information.

The M2002A and M2002B appliances are supplied with 4x40Gb Ethernet network links, 6x10Gb Ethernet network links and 8x1Gb Ethernet network links. If the appliances are configured for redundant HA, 2x1Gb links would be reserved for use by the appliance in addition to another interface to perform the HA replication (this can be configured to use any of the interfaces available or indeed an aggregated interface), leaving a potential total of 106Gb for customer workloads. In nonHA mode, all 148Gb connectivity can be utilised for workload traffic. There are a further two separate 1Gb links that are explicitly reserved for appliance administration. There are two modules that support 40Gb network connectivity with two ports available in each. There is a capacity limit of 40Gb per module. This report utilises 2 of the 40Gb links for workload traffic; for nonHA workload this is distributed over the two modules, for HA scenarios one module is utilised for workload traffic and one module for replication traffic.

All of the scenarios featured in this report utilise Request Responder messaging scenarios and the published messaging rate is measured in round trips/sec, which involves 2 message puts and 2 message gets. If you are only utilising one-way messaging (using a message sender, queue and message receiver to perform 1 message put and 1 message get), and you can avoid queue-lock contention, then you may achieve up to double the published rates.

The version of the MQ Appliance as tested in this report is M2002A MQ 9.2 and where a comparison is made to the restricted appliance configuration, this uses the MQ Appliance M2002B MQ 9.2.

## 4 Request/Responder Scenario

The scenario that will be used in this report reflects the most common usage patterns that customers are anticipated to use with the MQ appliance and provide guidance for those customers performing capacity planning or migration activities.

Each test was initially conducted and graphs produced using a 2K (2048 byte) message size. Additional tests were also conducted using 256byte, 20K and 200K to provide further data for capacity planning and are found in the accompanying data table.

As customers replace their existing MQ QM infrastructure, they may consolidate their MQ configuration from separate MQ QM servers (possibly running on different hardware and different MQ Versions) onto a single MQ appliance. They may have a mix of applications tightly bound to their existing QM and also a number of applications that connect using the MQ client API. To migrate to the MQ appliance all applications will need to connect via the MQ client API.

The following tests use MQ client connections and present the performance of MQ as deployed on the Appliance.

The test scenario in Figure 1 is a Request Responder scenario that simulates up to ten applications that interact with a single QM. A request queue and a reply queue will be created for each application, so ten pairs of queues are created for this test. One or more requester applications will send messages to one of the application request queues and will wait for a reply on the associated reply queue. Responder applications will listen for messages on the request queues before sending them to the correct reply queue.

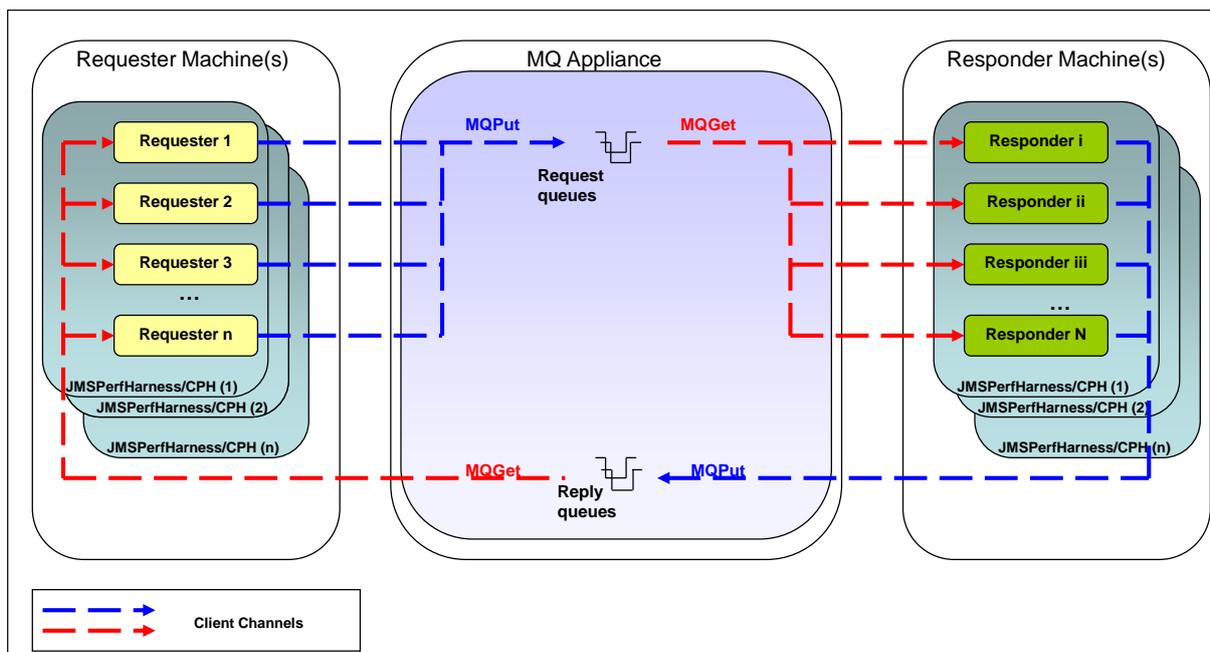


FIGURE 1 - REQUESTER-RESPONDER WITH REMOTE QUEUE MANAGER ON MQ APPLIANCE

Subsequent requester applications will send and receive messages from the set of application queues on a round-robin basis i.e. distributing the messages produced and consumed across the set of application queues.

Results are presented for various numbers of producer threads distributed across the 10 applications (using 10 pairs of queues), 200 fixed responder threads (20 responders per request queue) will send the replies to the appropriate reply queue, and the report will show the message rates achieved (in round trips/second) as the number of producers is increased.

#### 4.1 Test Scenario C1 – 10 Applications per QM, 1 QM, Non-persistent

The following graph shows how the scenario detailed in section 4 performs with Non-persistent messaging against a single QM.

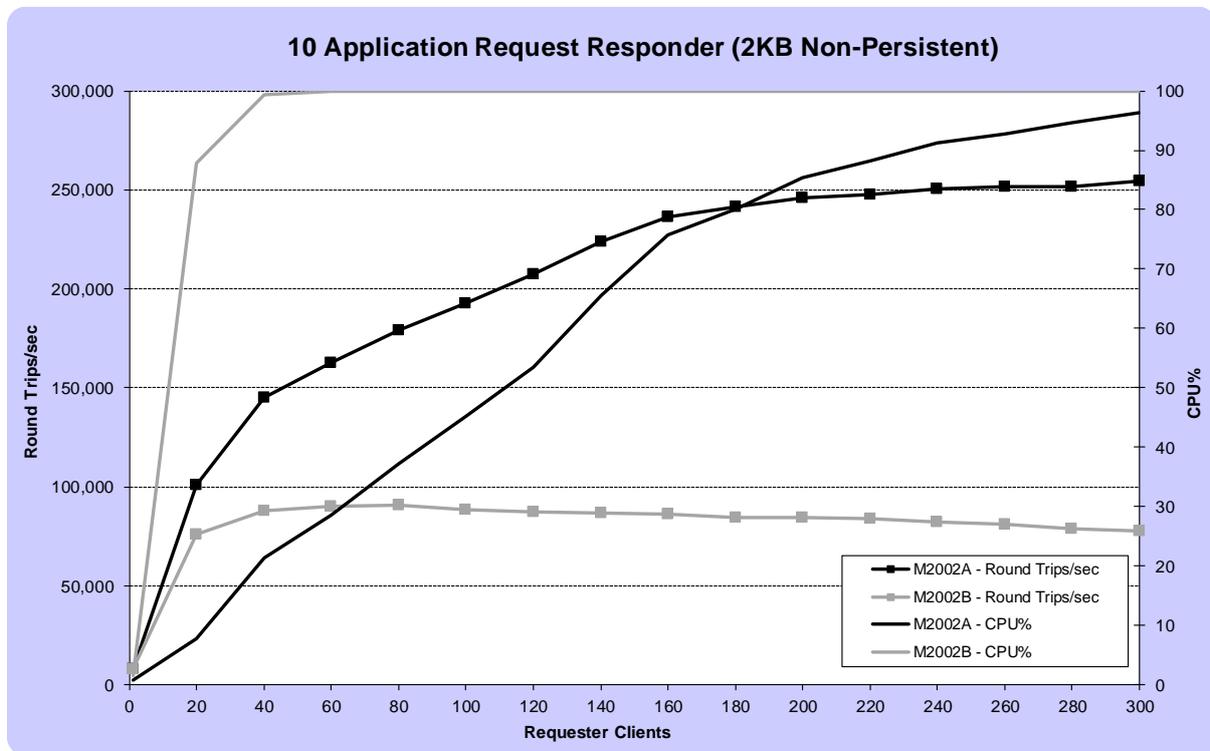


FIGURE 2 – PERFORMANCE RESULTS FOR 2KB NON-PERSISTENT MESSAGING

Figure 2 shows how by increasing the workload on the appliance (by increasing the number of concurrent requester clients), the throughput rate increases until the CPU capacity of the appliance is exhausted.

The M2002A appliance can achieve over 250,000 Round trips/sec for 2KB message size.

Test	M2002A			M2002B		
	Max Rate*	CPU%	Clients	Max Rate*	CPU%	Clients
10Q Request Responder (256b Non-persistent)	286,279	88.82	240	109,784	99.9	60
10Q Request Responder (2KB Non-persistent)	254,390	96.3	300	90,955	99.97	80
10Q Request Responder (20KB Non-persistent)	169,534	98.07	200	60,100	99.93	60
10Q Request Responder (200KB Non-persistent)	22,219	67.69	70	10,559	93.66	20

\*Round trips/sec

TABLE 1 - PEAK RATES FOR NON-PERSISTENT MESSAGING

## 4.2 Test Scenario C2 – 10 applications per QM, 1 QM, Persistent

This test repeats the test C1 featured in section 4.1, but utilises persistent messaging on the appliances local RAID10 disk subsystem.

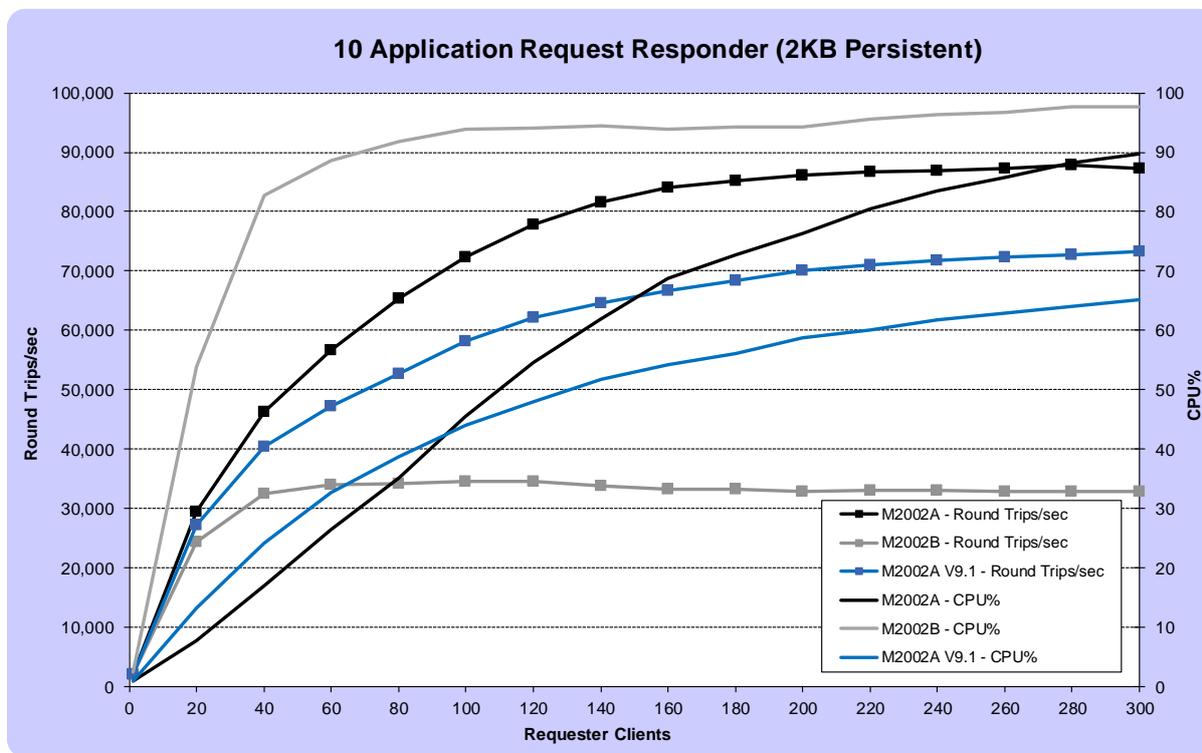


FIGURE 3 - PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING

Figure 3 shows that as the workload increases, a maximum throughput is achieved (over 87,000 Round trips/sec for 2KB message size) and the limits of the local disk subsystem have become the limiting factor for a single QM. This is an improvement of up to 20% when compared with MQ 9.1.

Test	M2002A			M2002B		
	Max Rate*	CPU%	Clients	Max Rate*	CPU%	Clients
10Q Request Responder (256b Persistent)	95,893	92.2	280	38,823	96.68	140
10Q Request Responder (2KB Persistent)	87,777	88.12	280	34,635	93.82	100
10Q Request Responder (20KB Persistent)	33,872	25.81	160	22,179	79.38	100
10Q Request Responder (200KB Persistent)	3,654	8.4	20	3,666	41.27	25

\*Round trips/sec

TABLE 2 – PEAK RATES FOR PERSISTENT MESSAGING

### 4.3 Test Scenario C3 – 10 applications per QM, 10 QM, Non-persistent

This test is equivalent to test C1 featured in section 4.1 with 10QM instead of 1QM. A total of 100 applications will be distributed across the 10 QM. This test demonstrates that there are no adverse effects from managing separate QMs within a single appliance.

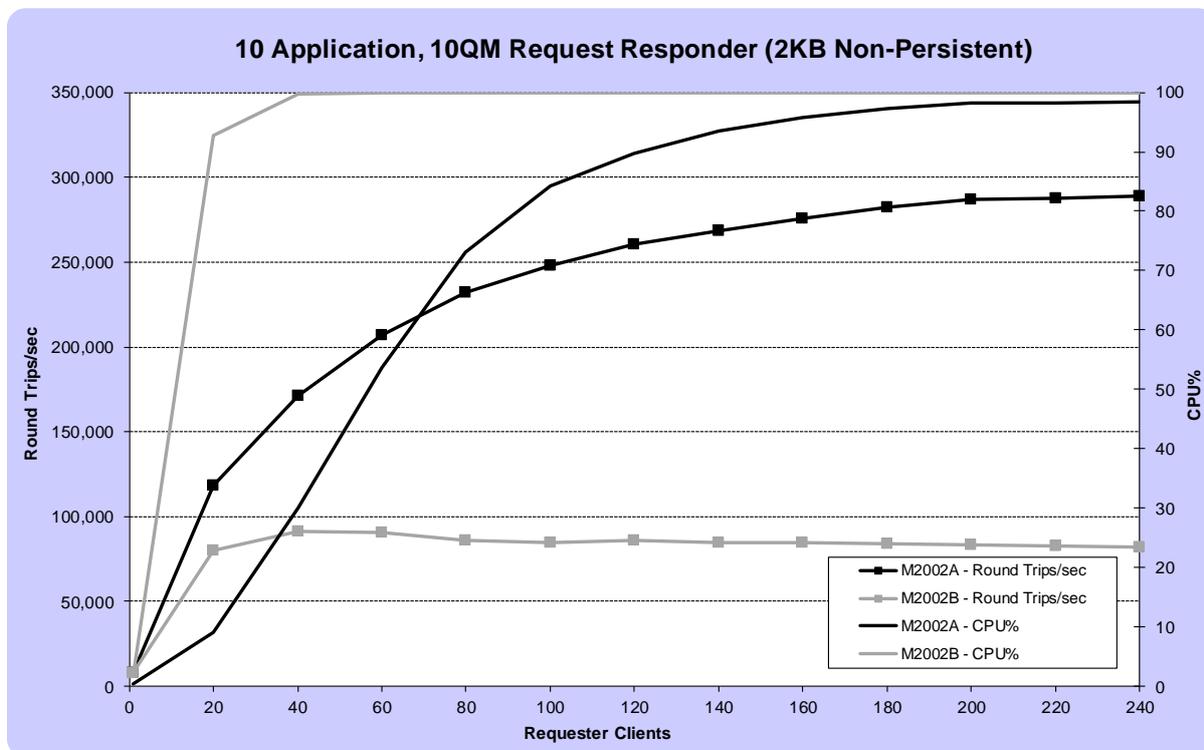


FIGURE 4 - PERFORMANCE RESULTS FOR 2KB, 10QM NON-PERSISTENT MESSAGING

Figure 4 shows that similar or improved performance can be obtained when running Non-persistent messaging through 10QM as compared with the single QM scenario.

The M2002A appliance can achieve approximately 3 times the throughput of the M2002B appliance.

Test	M2002A			M2002B		
	Max Rate*	CPU%	Clients	Max Rate*	CPU%	Clients
10Q 10QM Request Responder (256b Non-persistent)	336,386	97.95	200	109,199	99.98	60
10Q 10QM Request Responder (2KB Non-persistent)	288,913	98.36	240	91,081	99.66	40
10Q 10QM Request Responder (20KB Non-persistent)	186,776	98.68	160	61,923	99.57	40
10Q 10QM Request Responder (200KB Non-persistent)	22,554	42.18	60	12,373	98.28	30

\*Round trips/sec

TABLE 3 - PEAK RATES FOR 10QM NON-PERSISTENT MESSAGING

#### 4.4 Test Scenario C4 – 10 applications per QM, 10 QM, Persistent

This test repeats the test C3 featured in section 4.3, but utilises persistent messaging on the appliances local RAID10 disk subsystem. The graph and the accompanying data table illustrate that to utilise all of the available IO capacity on the appliance, multiple QM are required.

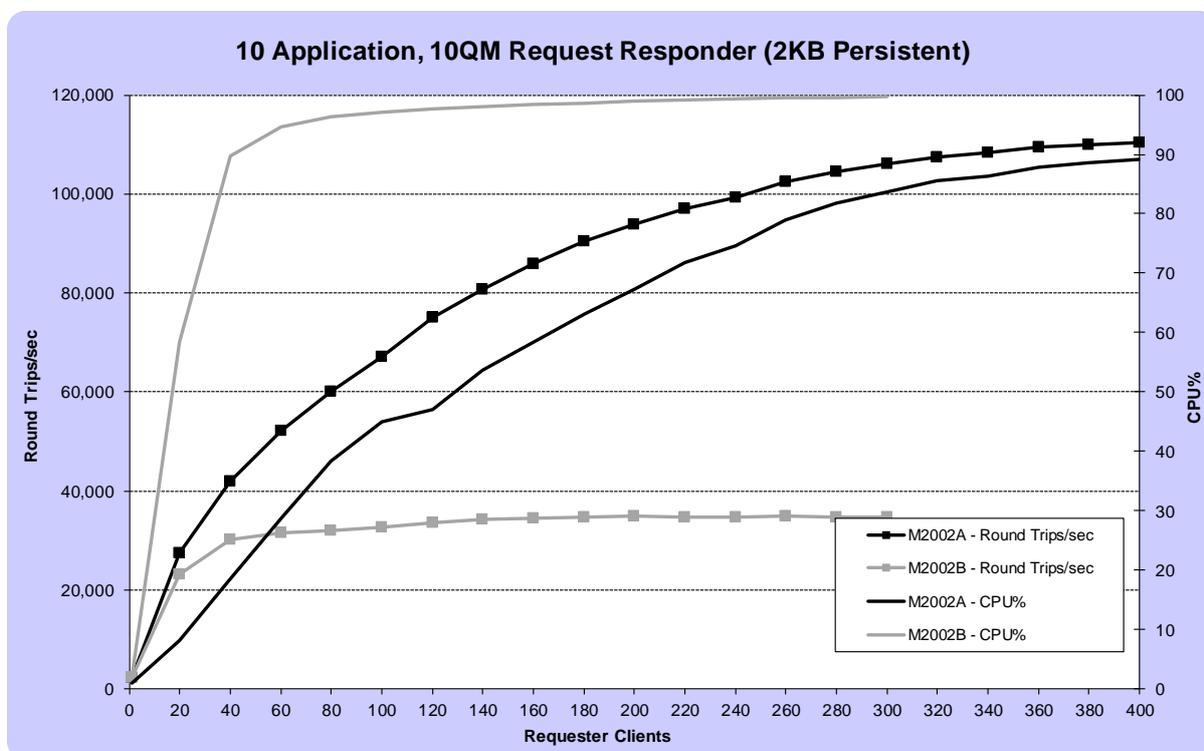


FIGURE 5 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING

Figure 5 shows that when we have multiple QM performing persistent messaging, the peak messaging rate obtainable on the M2002A is over 110,000 Round trips/sec for 2KB message size.

If using a message size of 2KB, the M2002A appliance can achieve over 3 times the persistent throughput of the M2002B appliance.

Test	M2002A			M2002B		
	Max Rate*	CPU%	Clients	Max Rate*	CPU%	Clients
10Q 10QM Request Responder (256b Persistent)	126,295	92.44	400	38,281	99.26	220
10Q 10QM Request Responder (2KB Persistent)	110,464	89.25	400	34,883	98.94	200
10Q 10QM Request Responder (20KB Persistent)	34,230	22.14	60	23,439	93.49	120
10Q 10QM Request Responder (200KB Persistent)	3,665	9.72	50	3,665	30.16	10

\*Round trips/sec

TABLE 4- PEAK RATES FOR 10QM PERSISTENT MESSAGING

## 5 Connection Scaling

The scaling measurement in this section is designed to test a scenario where there are a larger number of clients attached. Whereas the previous tests are optimised for throughput, these tests define an operational environment or scaling challenge to test from a performance perspective.

### 5.1 Connection Test

This test uses the Requester Responder workload as described in section 4. The requester applications are rated at 1 message every 100 seconds and 60,000 client bound requester applications are connected as fast as possible to determine the overall connection time for those clients to the MQ Appliance.

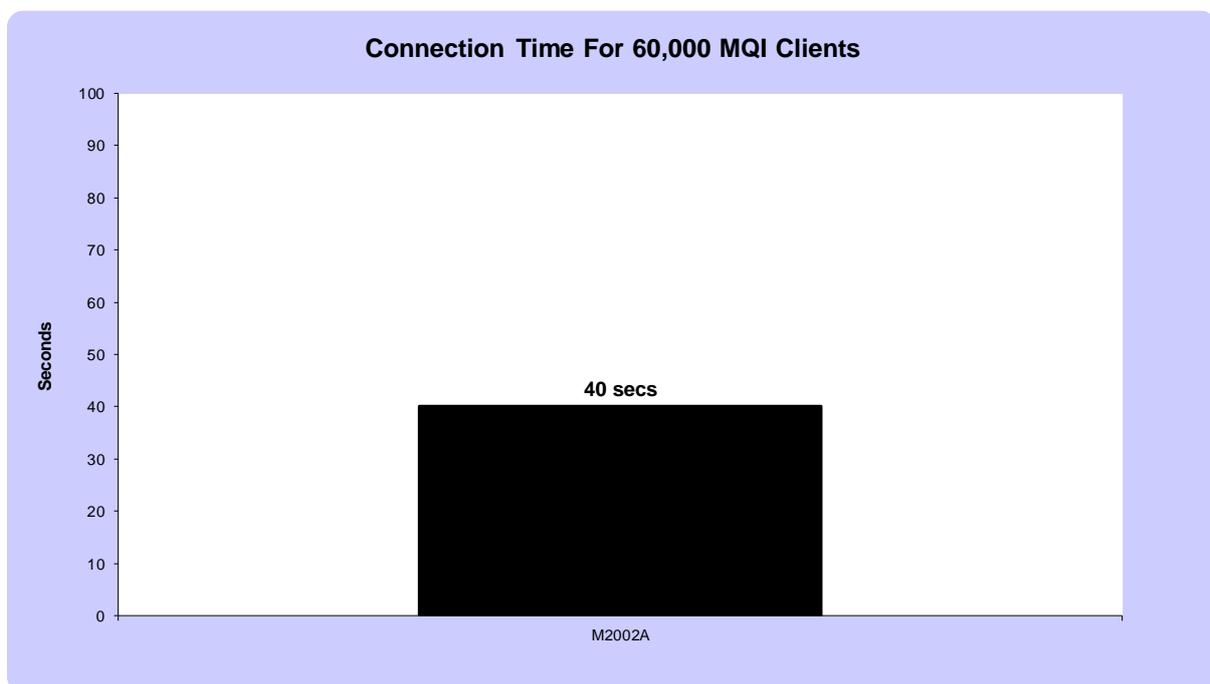


FIGURE 6 – PERFORMANCE RESULTS FOR MQI CLIENT CONNECTION TEST

The newer hardware of the M2002 appliance and the availability of 40Gb networking have made further significant improvements with regard to the rate at which the queue manager can accept connections from client bound applications and the MQ appliance can now support 60,000 clients initiating a connection to a single queue manager in approximately 40 seconds.

## 6 HA Scenarios

High Availability (HA) can be enabled by pairing two MQ Appliances together to provide continuous availability in the event of one of the appliances suffers a failure. The Queue Manager (QM) log and queue files are synchronously replicated across the pair of appliances.

If separate networks (and switches) are used to connect the pair of appliances, then the pair can also continue to operate in the event of a partial network outage.

To ensure clients reconnect to the QM on either of the pair of appliances, the clients should be made aware of the IP addresses assigned to the workload interfaces of both appliances; or a Virtualised IP address in the case that a suitable load balancer component is employed; or a floating IP if it is configured on the appliance for the QM.

To illustrate the performance profile of enabling the HA infrastructure, tests will be performed on two of the scenarios featured earlier in the report.

- 1) Request Responder 1QM Persistent (Test C2)
- 2) Request Responder 10QM Persistent (Test C4)

Each test will be conducted with both a standalone QM and a QM incorporated into an appliance HA group, so that the cost of the synchronous replication can be evaluated.

This section utilises the following connections:

Primary Appliance	Secondary Appliance	Notes
eth13	eth13	Connected directly between appliances with 1Gb copper patch cable
eth17	eth17	Connected directly between appliances with 1Gb copper patch cable
eth20-eth23	eth20-eth23	Unused
eth31	eth31	Connected directly between appliances with 40Gb copper cable
eth32		Workload driven via this 40Gb interface
eth33		Workload driven via this 40Gb interface

## 6.1 Test Scenario HA1 – 10 Applications per QM, 1 QM, Persistent

This test is identical to test C2 in section 4.2 and is presented here with results from running tests against a standalone QM and also against a QM that is included in an HA group.

Results are presented for various numbers of requester threads distributed across the 10 applications (using 10 pairs of queues), 200 fixed responder threads (20 responders per request queue) will send the replies to the appropriate reply queue, and the report will show the message rates achieved (in round trips/second) as the numbers of requesters are increased.

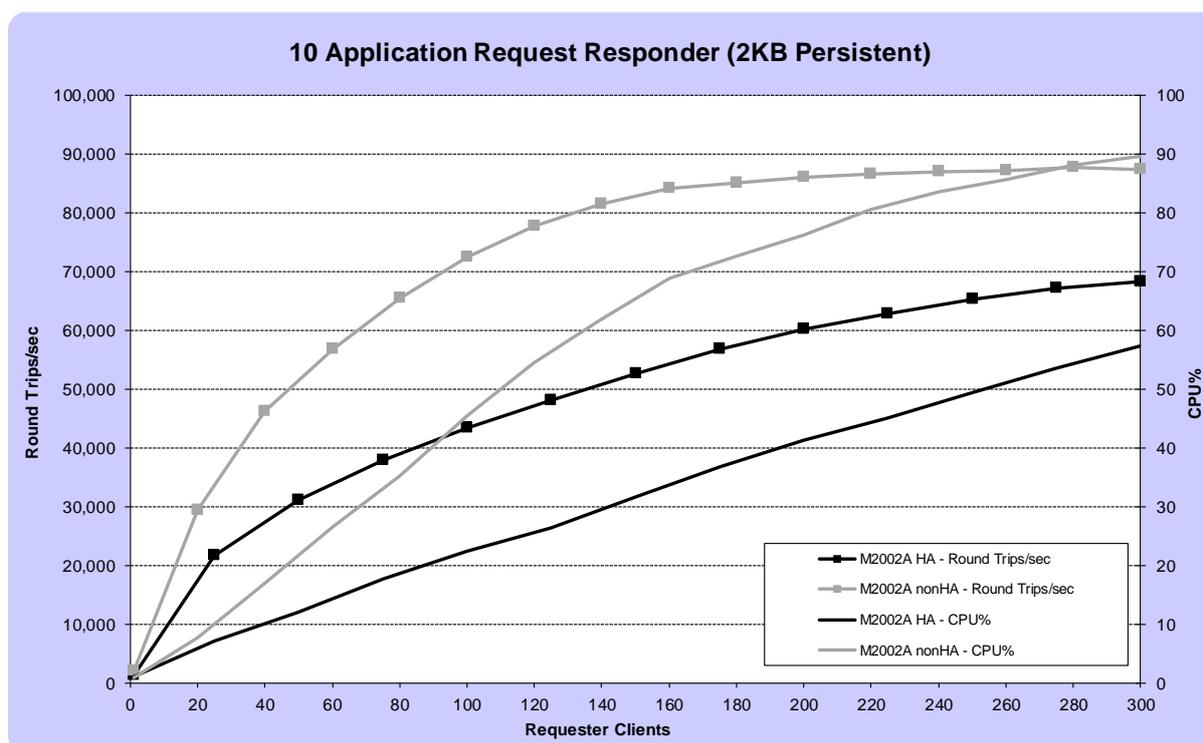


FIGURE 7 – PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING

Figure 7 shows that by enabling HA capability, the maximum throughput achieved in this single queue manager scenario with a 2K message size is reduced. There is a similar reduction in CPU utilisation, thus providing capacity to perform additional messaging on the appliance.

Like the scenario in section 4.2, this scenario has also increased and for the 2K message size, the performance of MQ V9.2 is up to 80% faster than MQ V9.1.

Test	M2002A HA				M2002A nonHA			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q Request Responder (256b Persistent)	86,947	83.57	400	0.7	95,893	92.2	280	0.4
10Q Request Responder (2KB Persistent)	72,461	67.53	400	0.8	87,777	88.12	280	0.5
10Q Request Responder (20KB Persistent)	16,753	13.36	140	0.9	33,872	25.81	160	0.6
10Q Request Responder (200KB Persistent)	1,750	4.49	15	2.0	3,654	8.4	20	1.4

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 5 - PEAK RATES FOR PERSISTENT MESSAGING

## 6.2 Test Scenario HA2 – 10 applications per QM, 10 QM, Persistent

This test repeats test C4 and is presented here with results from running tests against a standalone set of Queue Managers and also against a set of Queue Managers that are included in an HA group.

Results are presented for various numbers of requester threads distributed across the 10 Queue Managers who each host 10 pairs of queues (representing 10 applications per QM), 200 fixed responder threads (2 responders per request queue) will send the replies to the appropriate reply queue which are subsequently received by the originating requester threads, and the report will show the message rates achieved (in round trips/second) as the numbers of requesters are increased.

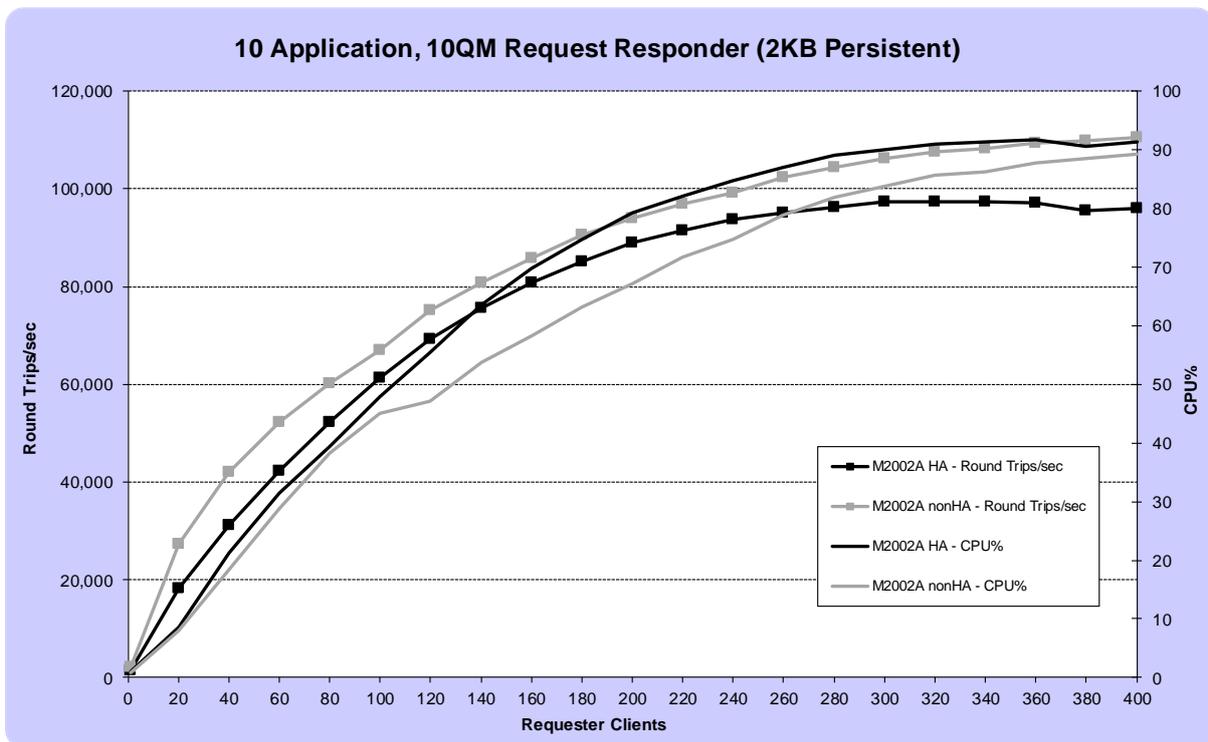


FIGURE 8 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING

Figure 8 shows that when we have multiple QM performing 2KB persistent messaging across a pair of HA appliances, the messaging rate is only 15% less than when distributed across a set of non HA Queue Managers.

Test	M2002A HA				M2002A nonHA			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q 10QM Request Responder (256b Persistent)	109,091	94.47	360	0.6	126,295	92.44	400	0.4
10Q 10QM Request Responder (2KB Persistent)	97,403	91	320	0.7	110,464	89.25	400	0.5
10Q 10QM Request Responder (20KB Persistent)	30,042	27.05	180	0.9	34,230	22.14	60	0.6
10Q 10QM Request Responder (200KB Persistent)	3,193	10.75	150	2.1	3,665	9.72	50	1.3

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 6 - PEAK RATES FOR 10QM PERSISTENT MESSAGING

### 6.3 How does HA perform over larger distances?

The previous section shows how the MQ appliance HA capability might perform if both appliances were located in the same data centre (i.e. 3m distance between the appliances). How would the HA performance differ if the pair of appliances were located a larger distance apart? Due to testing limitations, we need to simulate the latency that might be experienced as the distances between the appliances grows.

If the appliances are located 100Km apart, you might expect the smallest increase in packet transmission latency for each leg to be calculated as follows:

$$\text{distance} / \text{speed} = \text{time}$$

$$100,000\text{m} / 300,000,000\text{m/s}^1 = 0.000333\text{s} = 333 \text{ microseconds}$$

There must also be an allowance for the refraction index of the cable

$$333 * 1.5 = 500 \text{ microseconds}$$

Switching hardware and non-linear cable routing will likely further increase the latency between the pair of HA appliances. It is currently advised to customers to site a pair of HA appliances so that the latency between the two appliances is no greater than 10ms and preferably within the same data centre.

A delay can be inserted into the sending network layer of both appliances to simulate such latency and let us examine how this impacts the HA performance. The following chart repeats test HA2 from section 6.2 and shows the effect of a 2ms round trip latency introduced into the network layer between the two HA appliances.

---

<sup>1</sup> Assuming speed of light to be  $3 \times 10^8 \text{m/s}$

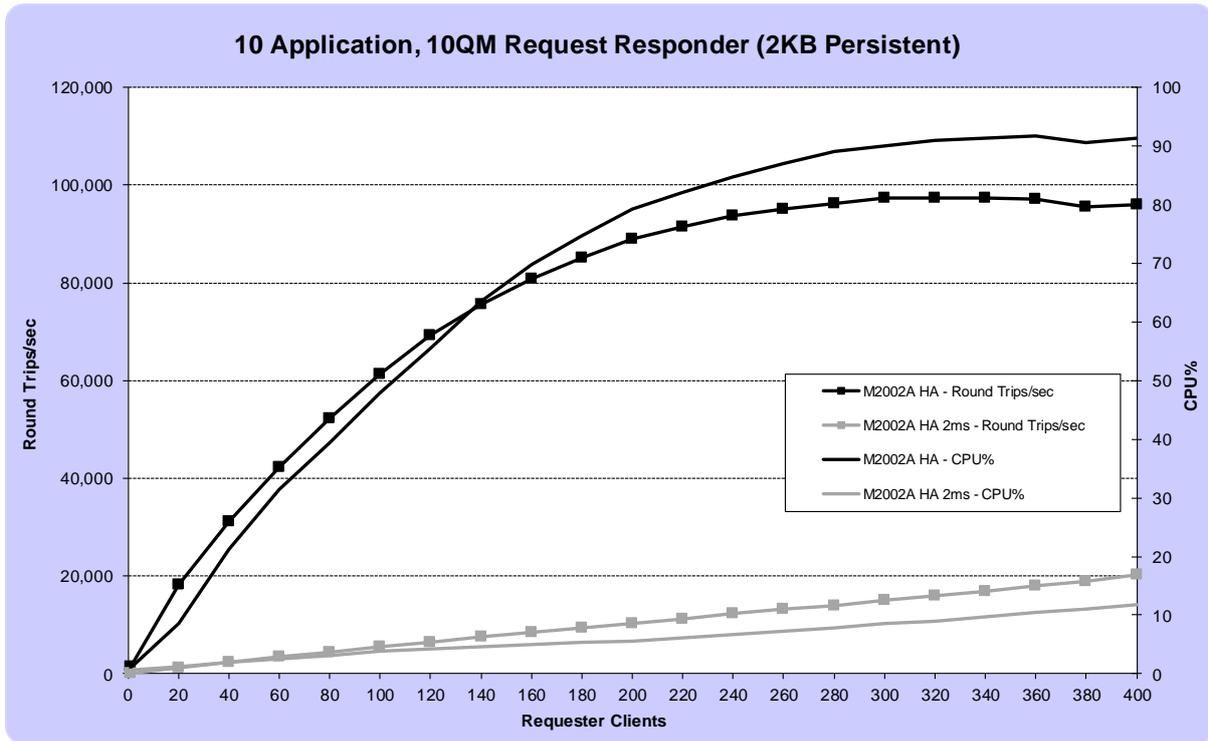


FIGURE 9 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING WITH/WITHOUT 2MS LATENCY

Figure 9 shows that an additional 2ms latency on the round trip time of the HA replication interface results in a ~80% reduction in performance than compared with the direct connection (no additional latency) between the appliances.

Test	M2002A HA				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 2ms Latency (256b Persistent)	20,331	10.62	400	8.8	18.6%
10Q 10QM Request Responder 2ms Latency (2KB Persistent)	20,285	11.76	400	10.5	20.8%
10Q 10QM Request Responder 2ms Latency (20KB Persistent)	17,727	17.1	400	11.2	59.0%
10Q 10QM Request Responder 2ms Latency (200KB Persistent)	3,038	12.16	150	12.4	95.1%

\*Round trips/sec  
#Single thread round trip latency (ms)

TABLE 7 - PEAK RATES FOR 10QM PERSISTENT MESSAGING WITH 2MS SIMULATED LATENCY

The data in the following tables show additional data points with simulated latency delays of 1, 5 and 10ms.

Test	M2002A HA				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 1ms Latency (256b Persistent)	37,724	21.46	400	4.7	34.6%
10Q 10QM Request Responder 1ms Latency (2KB Persistent)	36,590	24.68	400	5.7	37.6%
10Q 10QM Request Responder 1ms Latency (20KB Persistent)	24,307	26.07	400	6.0	80.9%
10Q 10QM Request Responder 1ms Latency (200KB Persistent)	3,199	14.49	120	7.1	100.2%

\*Round trips/sec  
#Single thread round trip latency (ms)

Test	M2002A HA				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 5ms Latency (256b Persistent)	8,684	5	400	20.8	8.0%
10Q 10QM Request Responder 5ms Latency (2KB Persistent)	8,400	5.46	400	25.2	8.6%
10Q 10QM Request Responder 5ms Latency (20KB Persistent)	7,262	8.06	400	26.6	24.2%
10Q 10QM Request Responder 5ms Latency (200KB Persistent)	1,623	6.08	150	27.6	50.8%

\*Round trips/sec

#Single thread round trip latency (ms)

Test	M2002A HA				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 10ms Latency (256b Persistent)	4,403	3.28	400	41.1	4.0%
10Q 10QM Request Responder 10ms Latency (2KB Persistent)	4,443	3.55	400	49.7	4.6%
10Q 10QM Request Responder 10ms Latency (20KB Persistent)	3,977	4.95	400	52.2	13.2%
10Q 10QM Request Responder 10ms Latency (200KB Persistent)	894	3.63	150	53.3	28.0%

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 8 - PEAK RATES FOR 10QM PERSISTENT MESSAGING WITH 1, 5 AND 10MS SIMULATED LATENCY

## 7 DR Scenarios

Users have the ability to configure a Queue Manager for Disaster Recovery (DR) to ensure that the QM data is distributed to a recovery appliance. This configuration allows the Queue Manager on the recovery appliance to resume work should an outage occur that results in the main appliance becoming unavailable.

Users have the ability to configure a QM for both HA and DR. The performance of this configuration will be examined in section 8, whilst in this section we will look at the standalone performance of DR.

The Queue Manager data is replicated asynchronously to the recovery appliance, which can result in messaging data loss (up to a maximum of 4MB per QM is held in the TCP send buffer) should the main appliance become unavailable. The Queue Manager at the recovery appliance must be manually started before it can start accepting connections from clients.

To illustrate the cost of enabling the DR infrastructure, tests will be performed on two of the scenarios featured earlier in this report.

- 1) Request Responder 1QM Persistent (Test C2)
- 2) Request Responder 10QM Persistent (Test C4)

Each test will be conducted with both a standalone QM and a QM configured with a remote DR appliance, so that the cost of the asynchronous replication can be evaluated.

This section utilises the following connections:

Appliance A	Appliance B	Appliance C (DR)	Notes
eth13	eth13		Connected directly between appliances with 1Gb copper patch cable. Used in section 8 only
eth17	eth17		Connected directly between appliances with 1Gb copper patch cable. Used in section 8 only
eth20-eth23	eth20-eth23		Unused
eth30		eth31	Connected directly between appliances with 40Gb copper cable for DR
eth31	eth31		Connected directly between HA appliances with 40Gb copper cable. Used in section 8 only
eth32			Workload driven via this interface
eth33			Workload driven via this interface

## 7.1 Test Scenario DR1 – 10 Applications per QM, 1 QM, Persistent

This test is identical to test C2 in section 4.2 and is presented here with results from running tests against a standalone QM and also against a QM that is configured for Disaster Recovery (although the recovery appliance is located 3m from the main appliance).

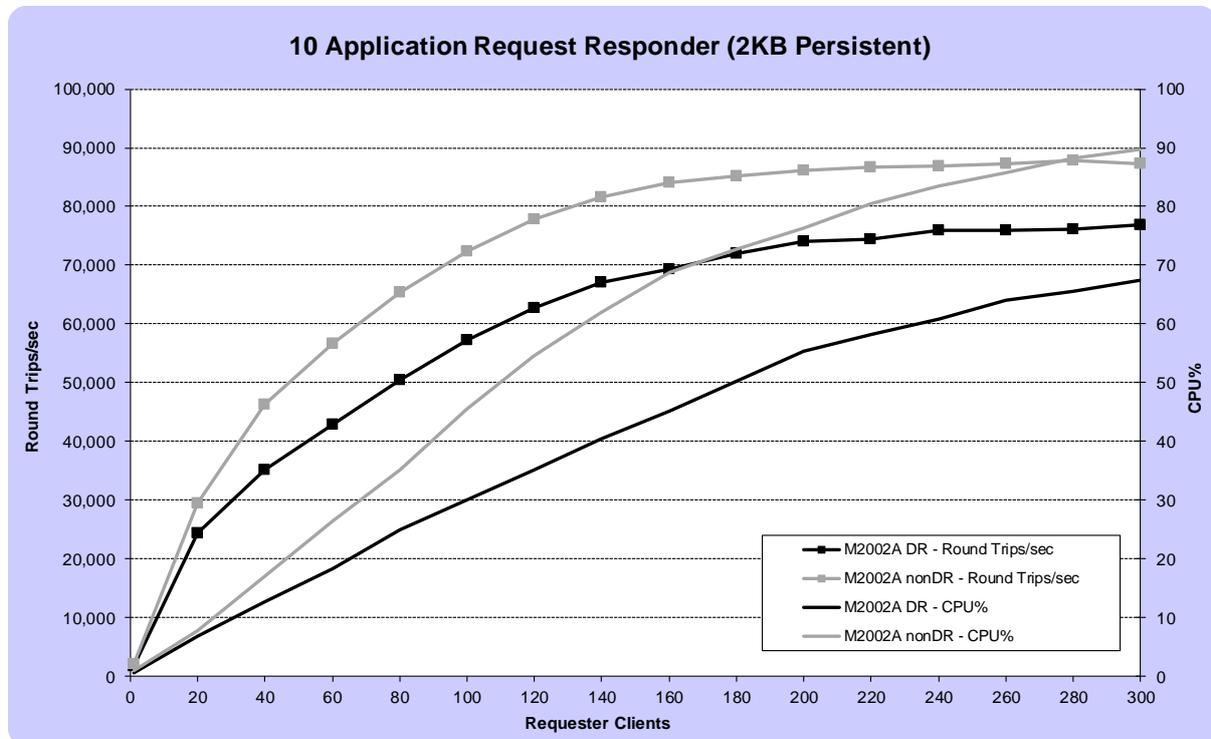


FIGURE 10 – PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING

Figure 10 shows that by enabling DR capability, the maximum throughput achieved with a 2K message size on a single Queue Manager is reduced by up to 15%. There is a similar reduction in CPU utilisation, thus providing capacity to perform additional messaging on the appliance.

Like the scenario in section 4.2, the throughput has also increased and for the 2K message size, the performance of MQ V9.2 is up to 35% faster than MQ V9.1.

Test	M2002A DR				M2002A nonDR			
	Max Rate*	CPU% Clients	Latency#		Max Rate*	CPU% Clients	Latency#	
10Q Request Responder (256b Persistent)	90,890	81.29	260	0.5	95,893	92.2	280	0.4
10Q Request Responder (2KB Persistent)	76,858	67.4	300	0.6	87,777	88.12	280	0.5
10Q Request Responder (20KB Persistent)	21,102	16.77	140	0.7	33,872	25.81	160	0.6
10Q Request Responder (200KB Persistent)	2,170	6.27	20	1.7	3,654	8.4	20	1.4

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 9 - PEAK RATES FOR PERSISTENT MESSAGING

## 7.2 Test Scenario DR2 – 10 Applications per QM, 10 QM, Persistent

This test is identical to the test C4 in section 4.4 and is presented here with results from running tests against ten standalone QM and also against ten QM that are configured for Disaster Recovery (although the recovery appliance is located 3m from the main appliance).

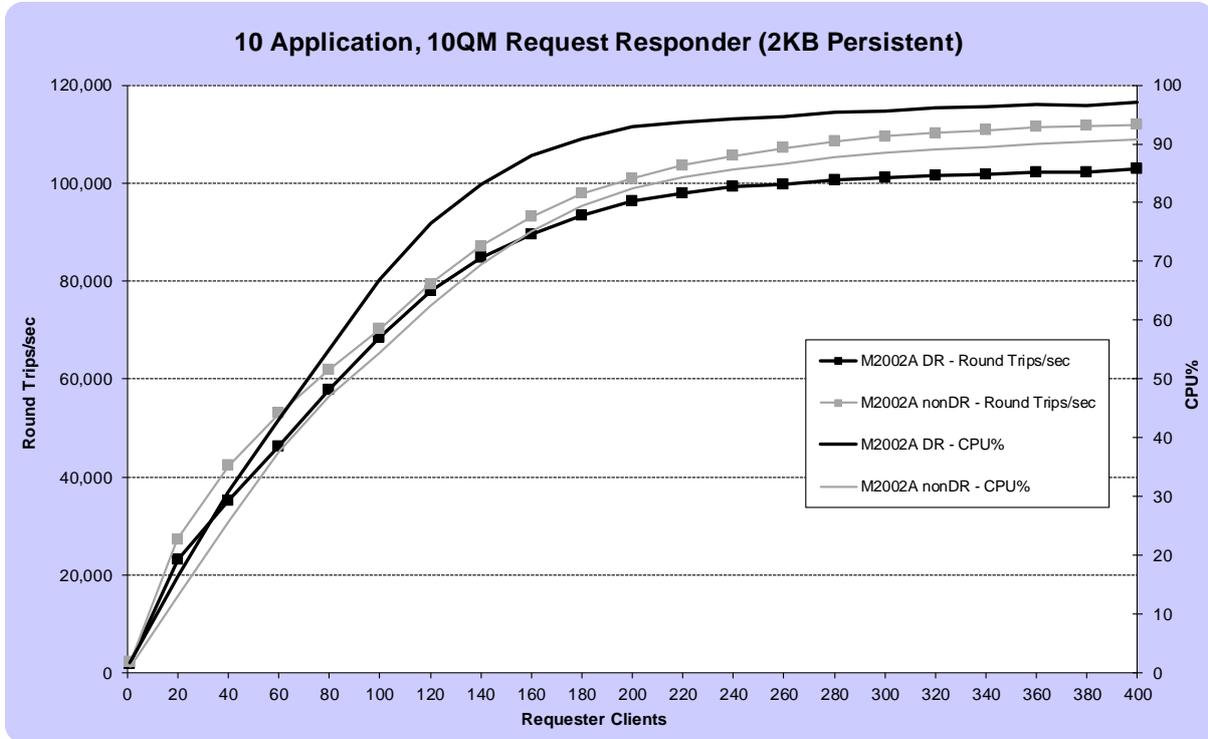


FIGURE 11 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING

Figure 11 shows that when we have multiple QM configured for Disaster Recovery performing 2KB persistent messaging, the peak messaging rate is within 10% of the rate achieved when distributed across a set of non DR Queue Managers.

Test	M2002A DR				M2002A nonDR			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q 10QM Request Responder (256b Persistent)	116,550	96.67	400	0.5	126,295	92.44	400	0.4
10Q 10QM Request Responder (2KB Persistent)	101,895	93.7	400	0.6	110,464	89.25	400	0.5
10Q 10QM Request Responder (20KB Persistent)	28,176	27.16	240	0.7	34,230	22.14	60	0.6
10Q 10QM Request Responder (200KB Persistent)	3,023	10.73	50	1.6	3,665	9.72	50	1.3

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 10 - PEAK RATES FOR 10QM PERSISTENT MESSAGING

### 7.3 How does DR perform over larger distances?

DR configuration usually requires the pair of appliances to be situated a large distance apart so that any particular event that might affect one appliance would be hoped not to affect the second appliance.

The data in the following tables show the results from the test scenario featured in the previous section but with additional data points using simulated latency delays of 10, 20, 50 and 100ms. A comparison against the DR scenario in which the MQ Appliances are directly connected is also included.

Test	M2002A DR				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 10ms Latency (256b Persistent)	113,738	95.68	360	0.5	97.6%
10Q 10QM Request Responder 10ms Latency (2KB Persistent)	99,083	88.33	380	0.6	97.2%
10Q 10QM Request Responder 10ms Latency (20KB Persistent)	21,346	16.84	240	0.7	75.8%
10Q 10QM Request Responder 10ms Latency (200KB Persistent)	2,162	6.69	50	1.8	71.5%

\*Round trips/sec

#Single thread round trip latency (ms)

Test	M2002A DR				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 20ms Latency (256b Persistent)	114,212	94.04	380	0.5	98.0%
10Q 10QM Request Responder 20ms Latency (2KB Persistent)	82,051	65.42	320	0.6	80.5%
10Q 10QM Request Responder 20ms Latency (20KB Persistent)	15,198	11.48	240	0.7	53.9%
10Q 10QM Request Responder 20ms Latency (200KB Persistent)	1,416	4.29	40	3.3	46.9%

\*Round trips/sec

#Single thread round trip latency (ms)

Test	M2002A DR				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 50ms Latency (256b Persistent)	98,510	75.68	360	0.5	84.5%
10Q 10QM Request Responder 50ms Latency (2KB Persistent)	54,290	36.82	400	0.7	53.3%
10Q 10QM Request Responder 50ms Latency (20KB Persistent)	9,155	7.87	300	1.3	32.5%
10Q 10QM Request Responder 50ms Latency (200KB Persistent)	898	3.05	50	7.9	29.7%

\*Round trips/sec

#Single thread round trip latency (ms)

Test	M2002A DR				vs Direct
	Max Rate*	CPU%	Clients	Latency#	
10Q 10QM Request Responder 100ms Latency (256b Persistent)	62,748	38.7	360	0.8	53.8%
10Q 10QM Request Responder 100ms Latency (2KB Persistent)	26,451	11.9	260	1.1	26.0%
10Q 10QM Request Responder 100ms Latency (20KB Persistent)	4,079	3.52	300	2.7	14.5%
10Q 10QM Request Responder 100ms Latency (200KB Persistent)	440	1.72	50	14.7	14.6%

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 11 - PEAK RATES FOR 10QM PERSISTENT MESSAGING WITH 10, 20, 50 AND 100MS SIMULATED LATENCY

## 8 HA and DR Scenarios

Users can configure a Queue Manager for both HA and DR. The data that is asynchronously replicated for disaster recovery is sent from the currently active instance of the HA pair.

This configuration allows the Queue Manager on the DR recovery appliance to resume work should an outage occur that results in both the appliances in the HA group becoming unavailable.

The performance of these scenarios is very close to that which has been measured in the HA scenarios featured in sections 6.1 and 6.2 showing that the cost of DR over and above the cost of HA is negligible (at least in direct connected scenarios).

Test	M2002A HA and DR				M2002A nonHA			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q Request Responder (256b Persistent)	84,082	83.92	450	0.7	95,893	92.2	280	0.4
10Q Request Responder (2KB Persistent)	69,305	65.56	400	0.8	87,777	88.12	280	0.5
10Q Request Responder (20KB Persistent)	15,393	13.58	120	0.9	33,872	25.81	160	0.6
10Q Request Responder (200KB Persistent)	1,625	5.17	15	2.1	3,654	8.4	20	1.4

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 12 - PEAK RATES FOR PERSISTENT MESSAGING, HA AND DR

Test	M2002A HA and DR				M2002A nonHA			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q 10QM Request Responder (256b Persistent)	111,944	92.09	400	0.7	126,295	92.44	400	0.4
10Q 10QM Request Responder (2KB Persistent)	98,636	90.3	400	0.8	110,464	89.25	400	0.5
10Q 10QM Request Responder (20KB Persistent)	29,627	32.31	180	0.9	34,230	22.14	60	0.6
10Q 10QM Request Responder (200KB Persistent)	3,148	12.78	40	2.0	3,665	9.72	50	1.3

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 13 - PEAK RATES FOR 10QM PERSISTENT MESSAGING, HA AND DR

## 9 Additional M2002A vs M2002B scenarios

In the earlier sections of this report detailing NonHA and NonDR scenarios, graphical and numerical data was also provided for the M2002B model. In the later sections detailing HA and DR scenarios, the comparison points were the equivalent NonHA and NonDR measurements.

This section has been added to illustrate the performance comparison of running HA, DR and HA+DR scenarios on either the M2002A or M2002B appliances.

We will initially look at the results of running test HA1 from section 6.1:

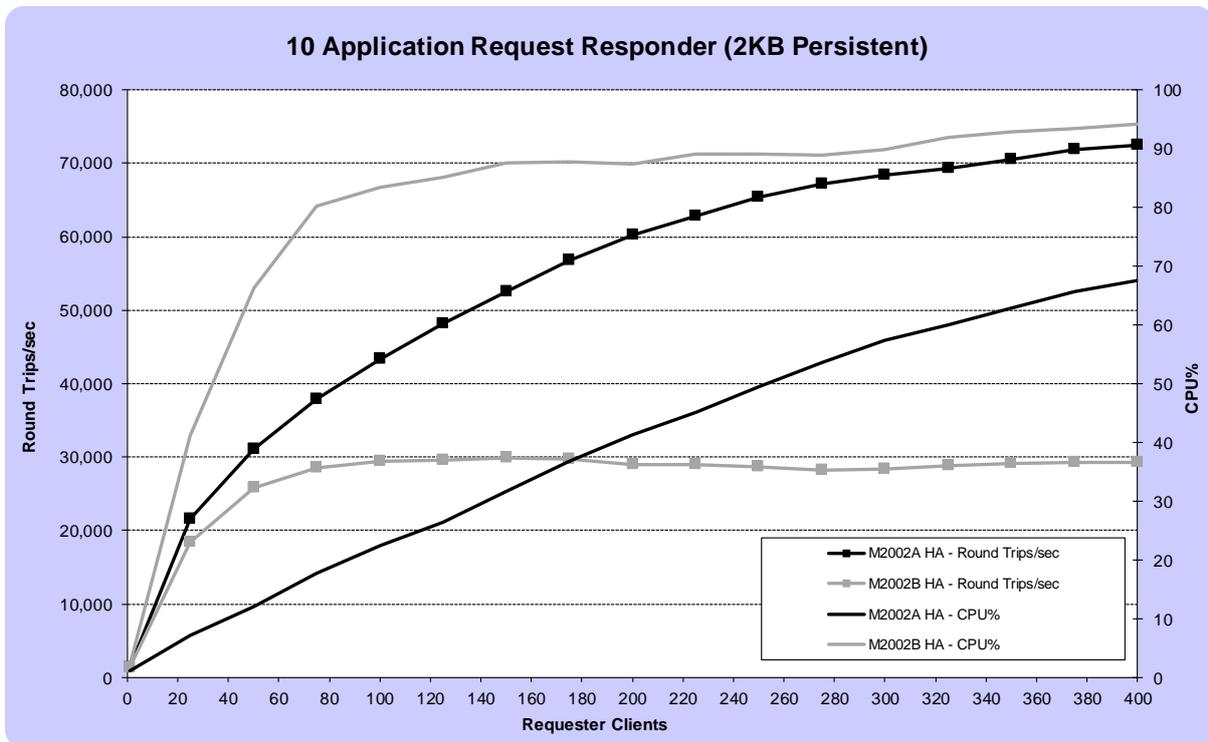


FIGURE 12 - PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING

Figure 12 shows that the M2002A appliance can achieve over twice the throughput of the M2002B appliance in this single HA QM scenario.

Test	M2002A HA				M2002B HA			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q Request Responder (256b Persistent)	86,947	83.57	400	0.7	33,346	92.02	150	0.7
10Q Request Responder (2KB Persistent)	72,461	67.53	400	0.8	29,891	87.63	150	0.8
10Q Request Responder (20KB Persistent)	16,753	13.36	140	0.9	15,693	58.15	120	0.9
10Q Request Responder (200KB Persistent)	1,750	4.49	15	2.0	1,966	22.37	20	2.0

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 14 - PEAK RATES FOR PERSISTENT MESSAGING

The following graph shows the results of running test HA2 from section 6.2:

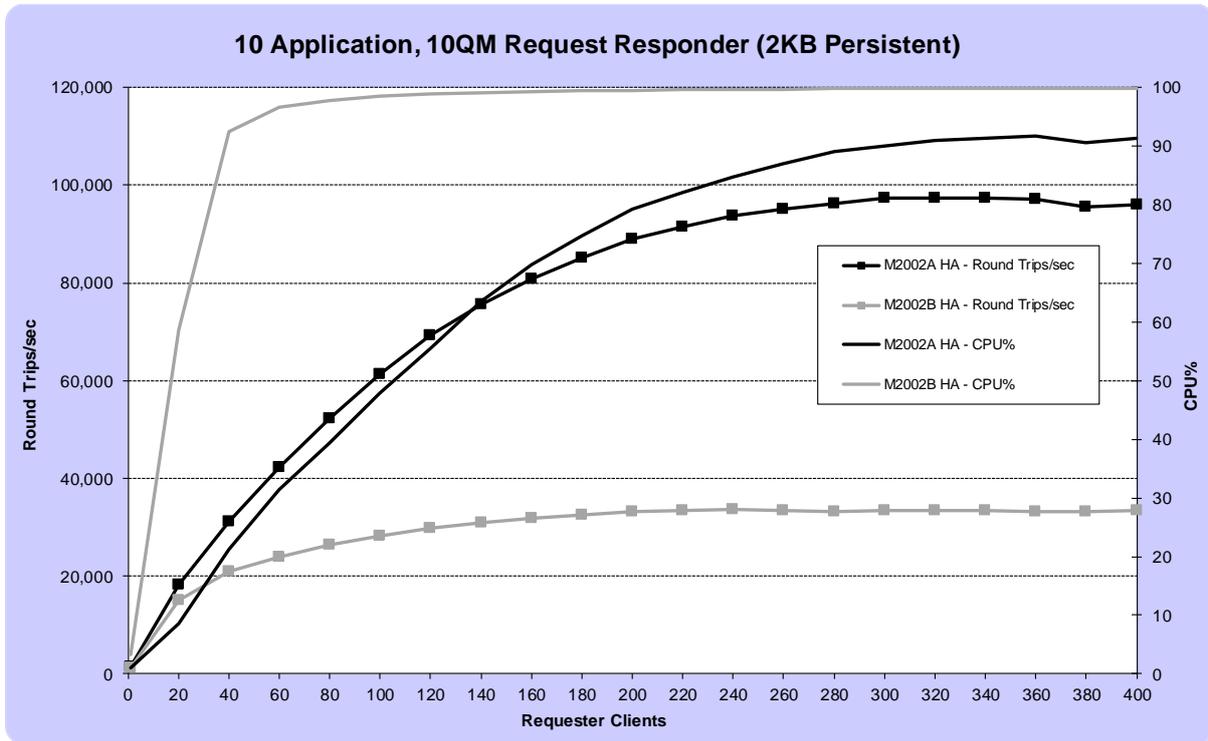


FIGURE 13 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING

Figure 13 shows that with multiple HA Queue Managers, the M2002A appliance can achieve approximately 3 times the throughput of the M2002B appliance.

Test	M2002A HA				M2002B HA			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q 10QM Request Responder (256b Persistent)	109,091	94.47	360	0.6	37,073	99.73	240	0.6
10Q 10QM Request Responder (2KB Persistent)	97,403	91	320	0.7	33,587	99.73	240	0.8
10Q 10QM Request Responder (20KB Persistent)	30,042	27.05	180	0.9	21,578	96.69	200	0.9
10Q 10QM Request Responder (200KB Persistent)	3,193	10.75	150	2.1	3,323	52.67	150	1.9

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 15 - PEAK RATES FOR 10QM PERSISTENT MESSAGING

The following graph shows the results of running test DR1 from section 7.1:

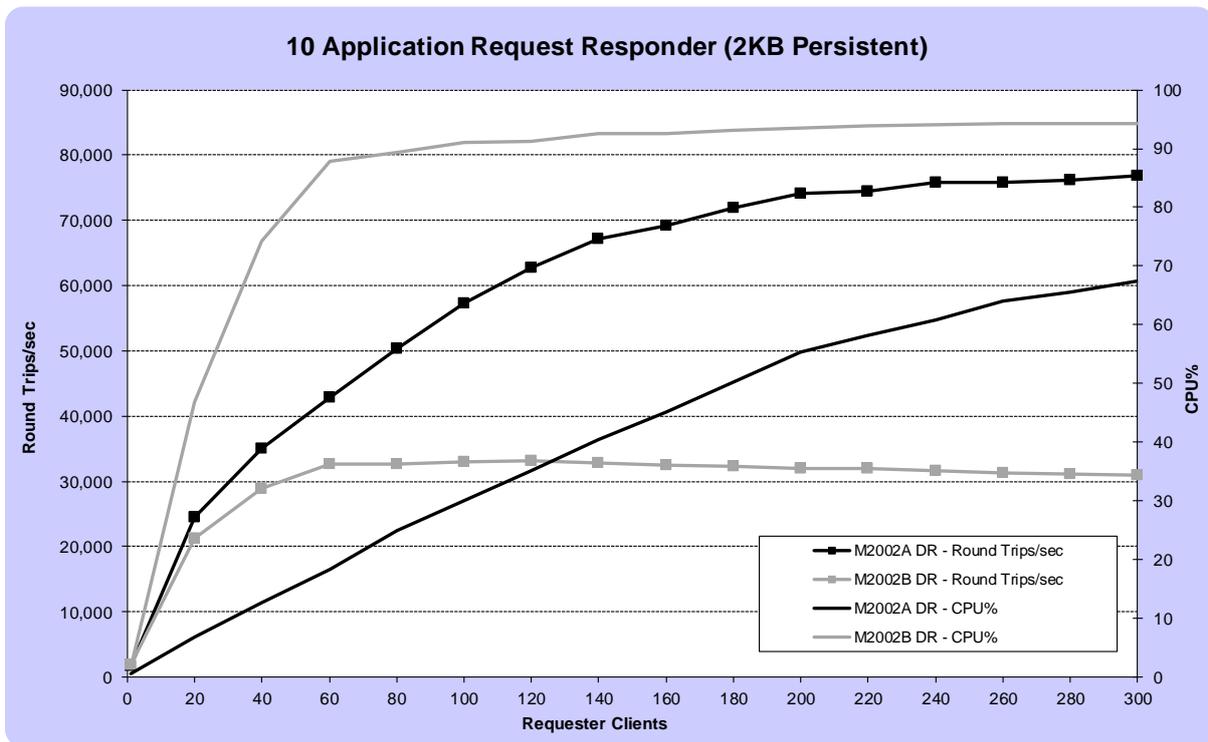


FIGURE 14 - PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING

Figure 14 shows that the M2002A appliance can achieve more than double the throughput of the M2002B appliance in this single DR QM scenario.

Test	M2002A DR				M2002B DR			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q Request Responder (256b Persistent)	90,890	81.29	260	0.5	37,925	93.92	120	0.5
10Q Request Responder (2KB Persistent)	76,858	67.4	300	0.6	33,172	91.2	120	0.5
10Q Request Responder (20KB Persistent)	21,102	16.77	140	0.7	15,716	61.82	100	0.7
10Q Request Responder (200KB Persistent)	2,170	6.27	20	1.7	2,097	31.18	20	1.6

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 16 - PEAK RATES FOR PERSISTENT MESSAGING

The following graph shows the results of running test DR2 from section 7.2:

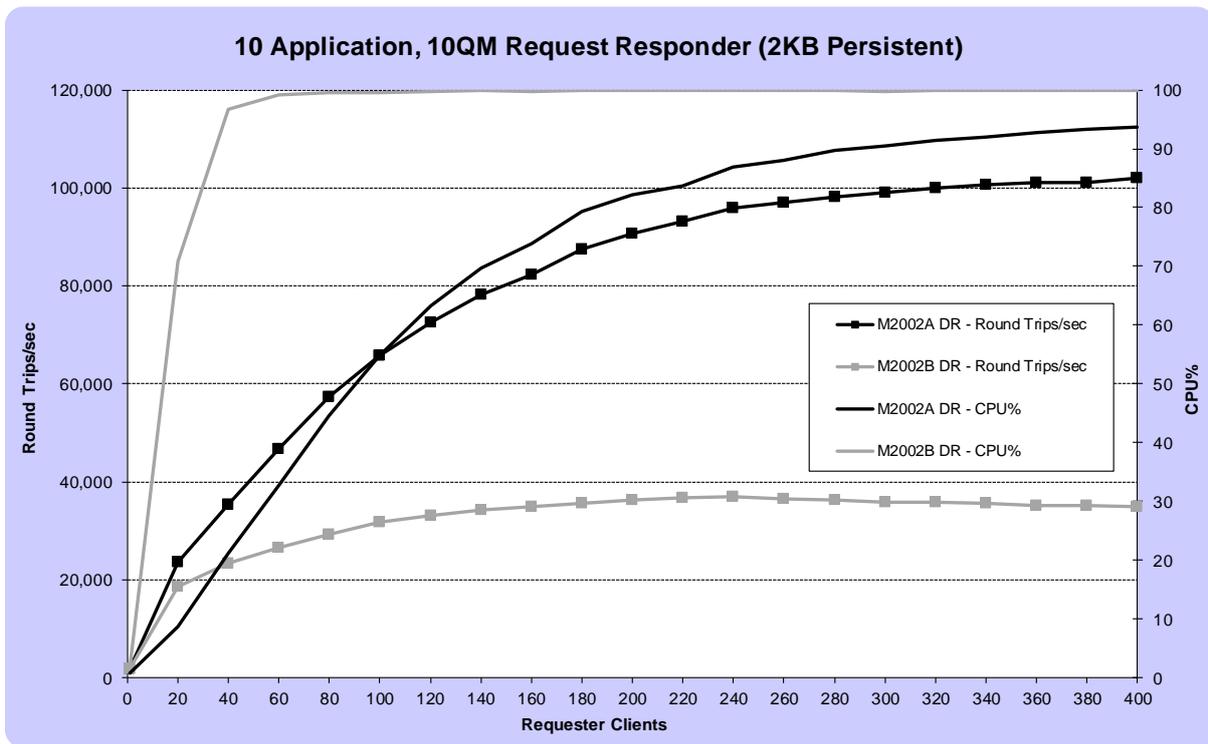


FIGURE 15 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING

Figure 15 shows that with multiple DR Queue Managers, the M2002A appliance can achieve nearly 3 times the throughput of the M2002B appliance.

Test	M2002A DR				M2002B DR			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q 10QM Request Responder (256b Persistent)	116,550	96.67	400	0.5	40,856	99.94	240	0.5
10Q 10QM Request Responder (2KB Persistent)	101,895	93.7	400	0.6	36,881	99.95	240	0.5
10Q 10QM Request Responder (20KB Persistent)	28,176	27.16	240	0.7	20,953	95.58	240	0.7
10Q 10QM Request Responder (200KB Persistent)	3,023	10.73	50	1.6	3,139	59.54	50	1.6

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 17 - PEAK RATES FOR 10QM PERSISTENT MESSAGING

The following graph shows the results of running the combined HA and DR scenario from section 8 for a single Queue Manager:

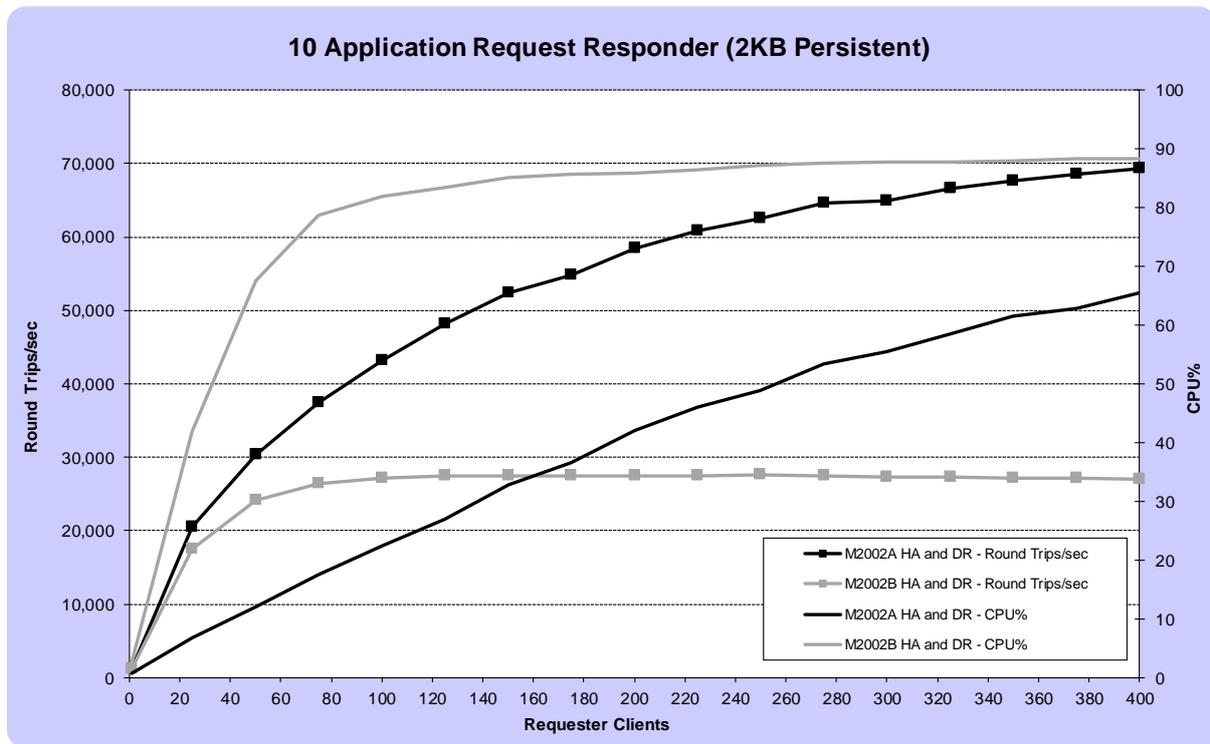


FIGURE 16 - PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING

Figure 16 shows that the M2002A appliance can achieve more than double the throughput of the M2002B appliance in this single HA and DR QM scenario.

Test	M2002A HA and DR				M2002B HA and DR			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q Request Responder (256b Persistent)	84,082	83.92	450	0.7	31,650	95.83	400	0.7
10Q Request Responder (2KB Persistent)	69,305	65.56	400	0.8	27,652	87.19	250	0.8
10Q Request Responder (20KB Persistent)	15,393	13.58	120	0.9	14,091	57.41	120	1.0
10Q Request Responder (200KB Persistent)	1,625	5.17	15	2.1	1,810	26.35	15	2.0

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 18 - PEAK RATES FOR PERSISTENT MESSAGING

The following graph shows the results of running the combined HA and DR scenario from section 8 against ten Queue Managers:

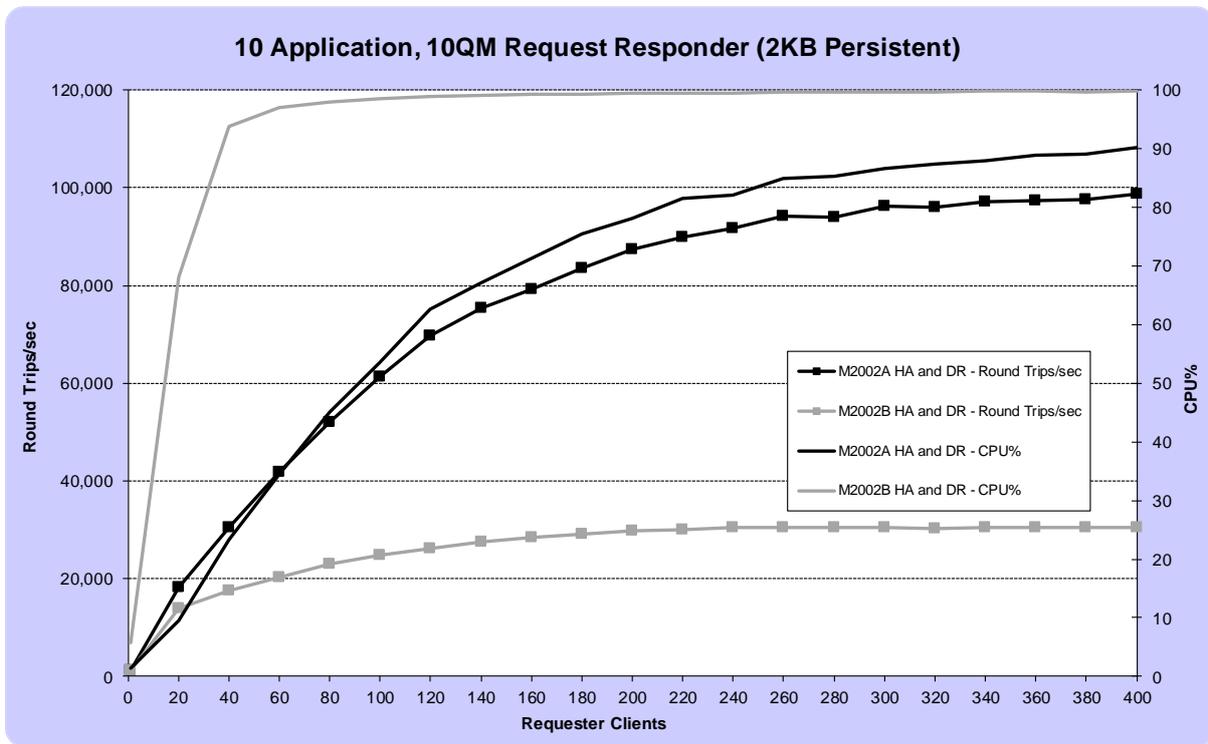


FIGURE 17 - PERFORMANCE RESULTS FOR 2KB, 10QM PERSISTENT MESSAGING

Figure 17 shows that with multiple HA and DR Queue Managers, the M2002A appliance can achieve over 3 times the throughput of the M2002B appliance.

Test	M2002A HA and DR				M2002B HA and DR			
	Max Rate*	CPU%	Clients	Latency#	Max Rate*	CPU%	Clients	Latency#
10Q 10QM Request Responder (256b Persistent)	111,944	92.09	400	0.7	35,461	99.6	300	0.7
10Q 10QM Request Responder (2KB Persistent)	98,636	90.3	400	0.8	30,472	99.73	300	0.8
10Q 10QM Request Responder (20KB Persistent)	29,627	32.31	180	0.9	19,164	97.05	400	0.9
10Q 10QM Request Responder (200KB Persistent)	3,148	12.78	40	2.0	3,136	66.75	150	2.0

\*Round trips/sec

#Single thread round trip latency (ms)

TABLE 19 - PEAK RATES FOR 10QM PERSISTENT MESSAGING

## 10 TLS

This section illustrates the cost of enabling TLS communication between the clients and the QM. We will use the scenario C1 from Section 4.1, and apply two of the strongest TLS 1.2 CipherSpecs to compare their performance.

Queue Manager authentication is used to setup the TLS conversation.

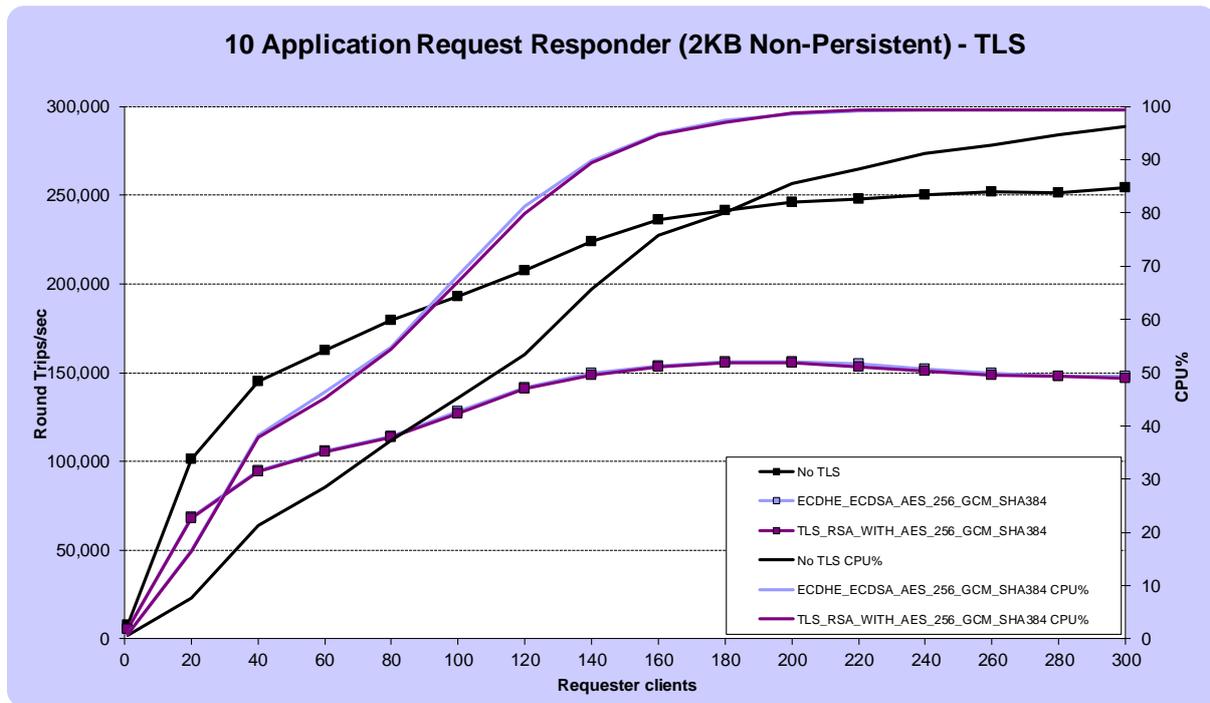


FIGURE 18 - PERFORMANCE RESULTS FOR 2KB NON-PERSISTENT MESSAGING WITH TLS

The graph in Figure 18 section illustrates that both the Elliptic Curve and RSA Ciphers using GCM (Galois/Counter Mode) perform equally at approximately 61% of the non-encrypted performance.

TLS Cipher	M2002A		
	Max Rate*	CPU%	Clients
ECDHE_ECDSA_AES_256_GCM_SHA384	156,055	98.58	200
TLS_RSA_WITH_AES_256_GCM_SHA384	155,833	98.69	200
No TLS	254,390	92.67	260
*Round trips/sec			

TABLE 20 - PEAK RATES FOR NON-PERSISTENT MESSAGING

## 11 AMS

This section illustrates the cost of enabling AMS to protect the message contents in transit between the clients and the QM and at rest at the QM. We will use the scenario C2 from Section 4.2, and compare Integrity, Privacy and Confidentiality mode with the Non-AMS performance.

The certificate key size used is 1024 bytes and the key reuse limit in Confidentiality mode was set to 32KB. The symmetric key encryption algorithm used was AES256. The cryptographic hash function used was SHA512.

The default certificate key size changed from 1024 to 2048 in MQ 9.1.4 and has an impact on AMS performance; the performance of Confidentiality mode with the increased key size is also shown.

Note that client CPU (rather than server CPU) is featured on the graph below as that shows the increase in computation performed by the clients in encrypting and decrypting the AMS protected messages.

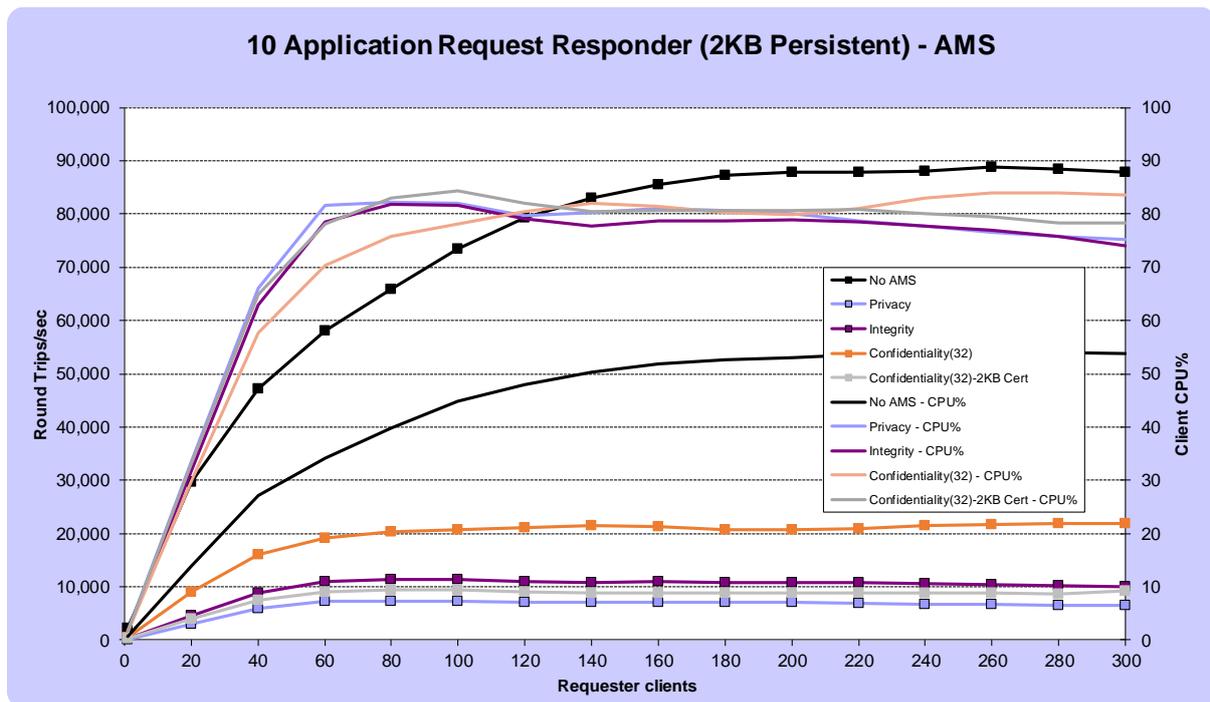


FIGURE 19 - PERFORMANCE RESULTS FOR 2KB PERSISTENT MESSAGING WITH AMS

The graph in Figure 19 shows that at 60 threads, the performance of AMS confidentiality is approximately 1/3 of the Non-AMS performance. The cost of the larger certificate key (resulting in more complex cryptographic calculations) further reduces the performance by half.

The peak throughput achieved for the AMS Confidentiality(32) measurement at 300 clients was nearly 22,000 round trip/s. The request/responder scenario utilizes a request and a reply queue, so for each round trip, 2 message put and 2 message get operations take place. For a single put/single get scenario, the peak performance that you might obtain in the same environment is over 44,000 msg/sec.

## 12 Frequently Asked Questions

### **Will I be able to use FASTPATH channels to send/receive messages into the MQ Appliance?**

Yes - this is now the default MQIBindType as specified in the Channels stanza in the qm.ini configuration file.

### **How do I view and change QM settings on the MQ Appliance?**

You can use the *dspmqini* command to view the QM configuration and *setmqini* to alter any configuration options. There are similar *dspmqvar* and *setmqvar* commands to view/alter environment variables.

### **What type of logging is used on the MQ Appliance?**

Only circular logging is supported on the MQ Appliance, and thus there are no facilities to monitor/prune QM logs.

### **Can I run my existing user exits?**

No – for appliance integrity, user exits will not be supported on the MQ Appliance. Many historic reasons for using code exits have now been resolved by product features.

### **What is throttling my messaging scenario?**

If customers experience throttled performance when driving high throughput workloads on M2002A, they should check the following:

- Persistent workloads - Customers might encounter the limits of the RAID10 subsystem as illustrated in this document
- Larger message (10K+) Non-persistent workloads - Customers might encounter network limits depending on which interfaces are selected for workload traffic. Customer can select higher bandwidth connectivity or aggregate multiple interfaces.
- Small message (2K-) Non-persistent workloads – Customers might encounter CPU saturation (Check MQ Console or CLI)

### **I have an M2001A/B, can I upgrade to M2002A/B?**

The M2002 appliance has completely different hardware to the M2001 appliance, therefore there is no upgrade option.

## 13 Appendix A – Client machine specification

The client machines (up to 4) used for the performance tests in this report have the following specification:

Category	Value
Machine	x3550 M5
OS	Red Hat Enterprise Linux Server 7.7
CPU	2x14 (2.6Ghz)
RAM	128GB RAM
Network	10/40Gb Ethernet
Disks	2x 120GB SAS SSD
RAID	ServeRAID M5210 (4GB Flash RAID cache) MQ Logs hosted on RAID-0 partition

## 14 Appendix B – QM Configuration

The following commands and expect scripts were used to create the standalone Queue Managers for this report:

```
crtmqm -lp 64 -lf 16384 -h 5000 -fs 16 PERFO
setmqini -m PERFO -s TuningParameters -k DefaultPQBufferSize -v 10485760
setmqini -m PERFO -s TuningParameters -k DefaultQBufferSize -v 10485760

proc configureQM { QMname QMport QMqueues } {
    send "runmqsc $QMname\n"
    send "define listener(L1) trrptype(tcp) port($QMport) control(qmgr)\n"
    send "start listener(L1)\n"
    send "alter channel(SYSTEM.DEF.SVRCONN) chltype(SVRCONN) sharecnv(1) maxmsgl(104857600)\n"
    send "alter qmgr maxmsgl(104857600)\n"
    send "alter qlocal(system.default.local.queue) maxmsgl(104857600)\n"
    send "alter qmodel(system.default.model.queue) maxmsgl(104857600)\n"
    send "alter qmodel(system.jms.model.queue) maxmsgl(104857600)\n"
    send "alter qmodel(system.jms.tempq.model) maxmsgl(104857600)\n"
    send "alter qlocal(system.dead.letter.queue) maxmsgl(104857600)\n"
    send "define channel(SYSTEM.ADMIN.SVRCONN) chltype(SVRCONN)\n"
    send "alter qmgr chlauth(disabled)\n"
    send "alter authinfo(SYSTEM.DEFAULT.AUTHINFO.IDPWOS) authtype(IDPWOS) chckclnt(OPTIONAL)\n"
    send "refresh security type(CONNAUTH)\n"
    send "define qlocal(queue) maxdepth(5000) replace\n"
    send "define qlocal(request) maxdepth(5000) replace\n"
    send "define qlocal(reply) maxdepth(5000) replace\n"
    for {set j 0} {$j <= $QMqueues} {incr j 1} {
        send "define qlocal(request$j) maxdepth(5000) replace\n"
        send "define qlocal(reply$j) maxdepth(5000) replace\n"
    }
    send "end\n"
}
```