

Cisco Nexus 9508: A New Speed Record in Data Center Switching

November 2016





TABLE OF CONTENTS

Executive Summary.....	3
About This Test	4
Performance Test Results	5
Test Bed Latency and Jitter	5
Testing Two Ports at a Time.....	7
Comparing Cisco Cloudscale Technology and Merchant Silicon ASICs.....	9
RFC 2889 Ethernet Unicast Performance	10
RFC 2544 IPv4 Unicast Performance	13
RFC 2544 IPv4 Unicast Performance With BGP Routing	15
RFC 5180 IPv6 Unicast Performance	17
RFC 5180 IPv6 Unicast Performance With BGP-MP Routing	19
RFC 3918 Ethernet Multicast Performance	19
Ethernet N-to-N Multicast Performance	21
RFC 3918 IPv4 Multicast Performance	27
IPv4 N-to-N Multicast Performance	29
Power Consumption.....	33
Forward Error Correction (FEC) Latency and Jitter	34
Test Methodology	36
Conclusion	38
Appendix A: Jitter Measurements.....	39
Appendix B: Software Releases Tested.....	50
About Network Test.....	50
Disclaimer.....	50



Executive Summary

Let's cut right to the chase: The Cisco Nexus 9508 is the densest, fastest data center switch we've ever tested. Even in the rarified world of data center core switches, Cisco's new Cloudscale Technology ASICs moved more traffic faster, and with lower latency and jitter, than any other switch we've evaluated.

For this large-scale (256-port) assessment of 100G Ethernet technology, Cisco Systems commissioned independent test lab Network Test to measure the performance of the Cisco Nexus 9508 switch. A key test component was the new N9K-X9732C-EX module, built around Cisco-designed Cloudscale Technology ASICs.

The Cisco Nexus 9508 proved its mettle across every high-stress test case. Among the major findings:

- Virtual line-rate performance in every test, regardless of frame size. In previous tests at this scale using merchant-silicon ASICs, the switch could not handle smaller frame sizes at line rate without loss
- Record low latency and jitter across all test cases
- Zero frame loss across all test cases covering unicast, multicast, Layer-2, Layer-3, and routing across all 256 100G Ethernet ports
- Zero frame loss in N-to-N multicast test cases, both in Layer-2 and Layer-3 configurations. In contrast, merchant silicon ASICs require gaps in test traffic to perform without loss
- 100G Ethernet power consumption as low as 23 watts per port
- Loss-free performance when forwarding to BGP and BGP-MP routes using IPv4 and IPv6
- Loss-free forwarding to more than 1 million IP multicast routes*OIFs/OILs in Layer-2 and Layer-3 configurations
- Wire-speed performance across all Layer-2 and Layer-3 multicast test cases

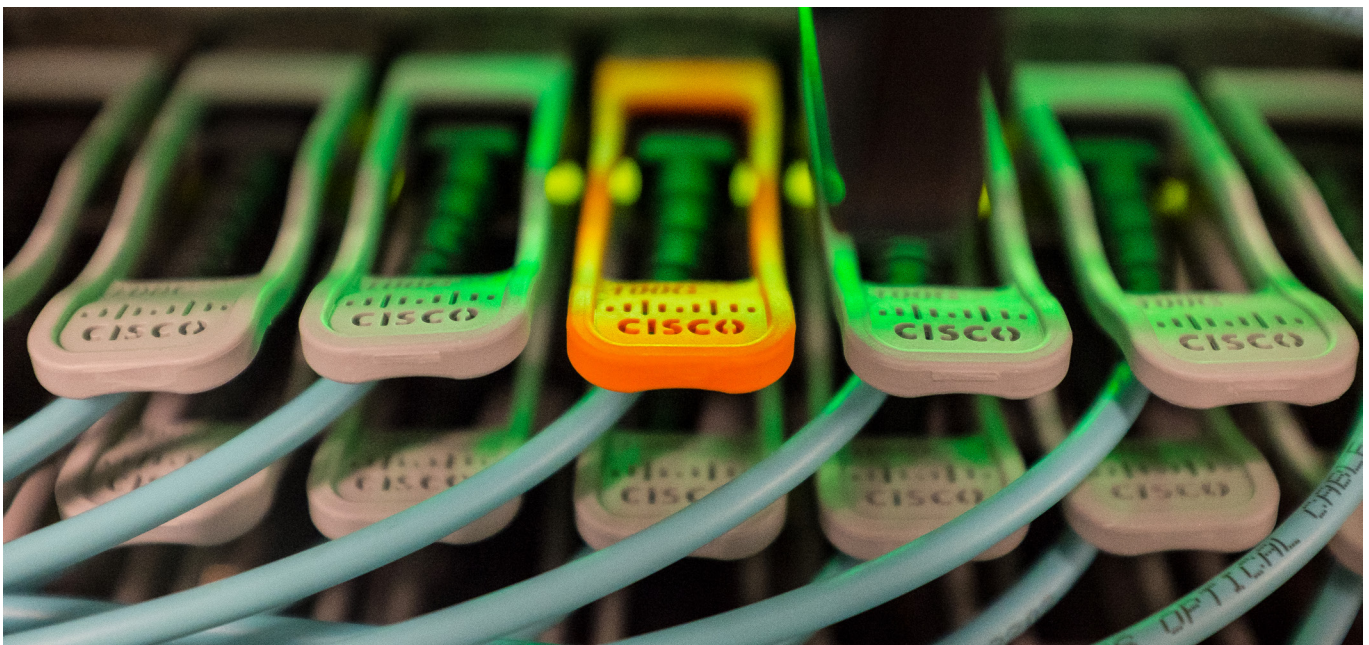


Figure 1: The Cisco Nexus 9508, with 256 100G Ethernet interfaces, on the test bed



About This Test

This project assessed the Cisco Nexus 9508 using 14 test cases, most involving 256 100G Ethernet interfaces:

- Test bed infrastructure latency and jitter
- Port-to-port performance
- Cisco Cloudscale Technology vs. merchant silicon throughput
- RFC 2889 Ethernet unicast performance
- RFC 2544 IPv4 unicast performance
- RFC 2544 IPv4 unicast performance with BGP routing
- RFC 5180 IPv6 unicast performance
- RFC 5180 IPv6 unicast performance with BGP-MP routing
- RFC 3918 Ethernet multicast performance
- Ethernet N-to-N multicast performance
- RFC 3918 IPv4 multicast performance
- IPv4 N-to-N Multicast Performance
- Power consumption
- Forward error correction (FEC) latency and jitter

The device under test for this project was a Cisco Nexus 9508 fully loaded with a Supervisor B engine and Cisco's new N9K-X9732C-EX modules, featuring Cisco Cloudscale Technology ASICs. The test bed, as seen in Figure 1, also included with the Spirent TestCenter traffic generator/analyzer with dX2-100G-P4 modules. The Spirent test instrument can offer traffic at wire speed on all ports with transmit timestamp resolution of 2.5 nanoseconds.

The primary metrics in this project were throughput, latency, and jitter.

[RFC 2544](#), the industry-standard methodology for network device performance testing, determines throughput as the limit of system performance. In the context of lab benchmarking, throughput describes the maximum rate at which a device forwards all traffic with zero frame loss.

Describing "real-world" performance is explicitly a non-goal of RFC 2544 throughput testing. Indeed, production networks load are typically far lower than the throughput rate.

Latency and jitter respectively describe the delay and delay variation introduced by a switch. Both are vital, and arguably even more important than throughput, especially for delay-sensitive applications such as video, voice, and some financial trading applications.

RFC 2544 requires latency be measured at, and only at, the throughput rate. Since average utilization in production networks is typically far lower than line rate, it can be useful to characterize delay for traffic at lower rates.

Accordingly, all tests described here present latency and jitter not only at the throughput rate, but also at 10, 50, and 90 percent of line rate. These results should help network professionals understand Cisco Nexus 9508 in *their* networks, modeling *their* network utilizations.

In all tests, engineers configured the Cisco Nexus 9508 with N9K-X9732C-EX modules in "cut-through" mode, its default setting. As described in the "Test Methodology" section, cut-through switching generally produces the lowest possible latency and jitter.



Performance Test Results

This section describes results for each configuration mode. See the “Test Methodology” section for details on test procedures.

Test Bed Latency and Jitter

Copper and fiber Ethernet transceivers and cabling add small but significant amounts of latency and jitter even when no switch is present. Cisco deployed a combination of CFP2 copper and fiber transceivers and adapters as well as 1- and 3-meter lengths of cabling – each of which adds delay outside the Cisco Nexus switch (each meter of copper or fiber cabling adds about 5 ns of delay).

Further, Cisco’s N9K-X9732C-EX modules use QSFP28 interfaces, while the Spirent TestCenter 100G Ethernet dX2 test modules use CFP2 interfaces. This difference necessitated the use of CFP2-to-QSFP28 electrical adapters, introducing significant additional delay.

To characterize the latency and jitter of these external factors, test engineers took measurements between ports of the Spirent TestCenter traffic generator/analyzer with no switch present. Engineers ran the tests twice, once apiece with copper and fiber media.

Tables 1 and 2 present results of test bed latency and jitter characterization for copper and fiber infrastructure (transceivers plus cabling), respectively. It’s important to note that the rest of the latency and jitter measurements in this report include this extra delay and delay variation. These measurements cannot easily be subtracted from switch measurements due to the use of fully meshed traffic (including fiber ports exchanging traffic with copper ports), the combination of copper and fiber transceivers, and different cable lengths. Still, these test-bed-only measurements offer guidelines for the additional latency and jitter present due to factors external to the Cisco Nexus 9508 switch.



Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	297,619,047.62	0.152	100.0000%	0.340	0.351	0.380	0.005	0.020
128	168,918,918.92	0.173	100.0000%	0.340	0.351	0.380	0.003	0.020
256	90,579,710.15	0.186	100.0000%	0.340	0.352	0.380	0.004	0.020
512	46,992,481.21	0.192	100.0000%	0.340	0.351	0.380	0.005	0.030
1,024	23,946,360.16	0.196	100.0000%	0.340	0.350	0.380	0.006	0.030
1,280	19,230,769.23	0.197	100.0000%	0.340	0.351	0.380	0.006	0.030
1,518	16,254,876.47	0.197	100.0000%	0.340	0.351	0.380	0.005	0.030
9,216	2,706,799.48	0.200	100.0000%	0.340	0.352	0.380	0.005	0.030

Table 1: Copper test bed infrastructure performance results

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	297,619,047.61	0.152	100.0000%	0.330	0.346	0.380	0.005	0.020
128	168,918,918.91	0.173	100.0000%	0.330	0.346	0.380	0.003	0.020
256	90,579,710.14	0.186	100.0000%	0.330	0.346	0.380	0.004	0.020
512	46,992,481.20	0.192	100.0000%	0.330	0.347	0.380	0.004	0.020
1,024	23,946,360.16	0.196	100.0000%	0.330	0.347	0.380	0.006	0.020
1,280	19,230,769.23	0.197	100.0000%	0.330	0.347	0.380	0.006	0.020
1,518	16,254,876.47	0.197	100.0000%	0.330	0.347	0.380	0.005	0.020
9,216	2,706,799.48	0.200	100.0000%	0.330	0.348	0.380	0.005	0.020

Table 2: Fiber test bed infrastructure performance results



Testing Two Ports at a Time

A good benchmark should be stressful, and accordingly most tests described here involved all 256 100G Ethernet ports using maximally stressful traffic patterns. Engineers also ran a few extra tests involving just two ports to characterize latency and jitter when the switch is less heavily loaded.

Engineers ran three sets of 2-port tests:

- Port to port on the same module and same ASIC¹
- Port to port on the same module with different ASICs
- Port to port across different modules

Tables 3, 4, and 5 present throughput, latency, and jitter results from the various 2-port test cases. Note that latency and jitter is very similar in same-module and cross-module ASIC-to-ASIC tests (see Tables 4 and 5). Expressed another way, there is no additional latency or jitter penalty for moving traffic between switch modules.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Gbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	297,601,203	152.372	99.994%	1.160	1.190	1.230	0.005	0.030
128	168,908,791	172.963	99.994%	1.160	1.192	1.240	0.003	0.020
256	90,574,279	185.496	99.994%	1.180	1.203	1.240	0.003	0.030
512	46,989,664	192.470	99.994%	1.180	1.207	1.240	0.004	0.030
1,024	23,944,924	196.157	99.994%	1.200	1.227	1.270	0.005	0.030
1,280	19,229,616	196.911	99.994%	1.200	1.227	1.270	0.006	0.030
1,518	16,253,902	197.387	99.994%	1.200	1.227	1.270	0.005	0.030
9,216	2,706,637	199.555	99.994%	1.190	1.227	1.270	0.005	0.040

Table 3: Same-ASIC performance results

¹ Each N9K-X9732C-EX module has 4 switch ASICs, or one ASIC per 8 front-panel ports.



Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Gbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	297,601,203	152.372	99.994%	2.640	2.676	2.720	0.005	0.030
128	168,908,791	172.963	99.994%	2.660	2.688	2.730	0.003	0.020
256	90,574,279	185.496	99.994%	2.690	2.713	2.750	0.003	0.030
512	46,989,664	192.470	99.994%	2.720	2.743	2.780	0.004	0.040
1,024	23,944,924	196.157	99.994%	2.770	2.799	2.850	0.004	0.040
1,280	19,229,616	196.911	99.994%	2.770	2.800	2.850	0.004	0.040
1,518	16,253,902	197.387	99.994%	2.770	2.799	2.840	0.004	0.030
9,216	2,706,637	199.555	99.994%	2.770	2.800	2.840	0.004	0.040

Table 4: ASIC-to-ASIC, same-module performance results

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Gbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	297,601,203	152.372	99.994%	2.640	2.689	2.740	0.005	0.030
128	168,908,791	172.963	99.994%	2.660	2.701	2.760	0.003	0.020
256	90,574,279	185.496	99.994%	2.690	2.725	2.780	0.003	0.030
512	46,989,664	192.470	99.994%	2.720	2.756	2.810	0.004	0.030
1,024	23,944,924	196.157	99.994%	2.770	2.812	2.870	0.004	0.040
1,280	19,229,616	196.911	99.994%	2.780	2.813	2.870	0.004	0.040
1,518	16,253,902	197.387	99.994%	2.770	2.811	2.870	0.004	0.030
9,216	2,706,637	199.555	99.994%	2.770	2.811	2.860	0.005	0.040

Table 5: ASIC-to-ASIC, cross-module performance results



Comparing Cisco Cloudscale Technology and Merchant Silicon ASICs

Cisco says the Cloudscale Technology ASICs it designed for its new N9K-X9732C-EX modules provide higher performance than earlier N9K-X9432C-S modules, which use merchant silicon ASICs. To validate that claim, Network Test compared the performance of the two systems using the same test instrument and test bed².

Figure 2 compares throughput for various test cases with minimal-length frames. The key difference is that the new Cisco modules deliver essentially line-rate throughput in all cases. Powered by Cisco-designed Cloudscale Technology ASICs, the N9K-X9732C-EX modules never dropped a frame in any test, regardless of frame size.

For purposes of comparison with earlier tests, test traffic carried not only IPv4 but also UDP headers. In earlier tests involving merchant-silicon ASICs, those switching chips employed a default hashing algorithm that used Layer-2/3/4 criteria to distribute flows across the switch fabric. In that case, engineers configured UDP headers to use 8,000 unique source and destination ports.

The new Cisco ASIC in the N9K-X9732C-EX hashes by default on Layer-2 and Layer-3 criteria, and thus does not require Layer-4 information for uniform flow distribution. Nonetheless, engineers retained the UDP headers to allow customers to make direct comparisons between the two types of ASICs.

Comparisons of latency and jitter are not possible because of different default configuration modes. The earlier modules use “store-and-forward” mode by default, meaning they cache each entire incoming frame before forwarding it³. In contrast, the default for the new modules is “cut-through” mode, meaning they begin forwarding each incoming frame as soon as it is received.

Because industry-standard test practices require different measurement methods for the two modes, direct comparisons are not meaningful. As a rule, cut-through mode generally provides lower latency and jitter, with the tradeoff that frames are not checked for errors. Since even large data center designs typically involve short distances between switches, and since the probability of data corruption increases with cable length, the risk of data corruption is relatively low.

The new modules also have significantly larger output buffers, as discussed in sections on Ethernet and IPv4 N-to-N multicast performance. In both test cases, the new module can accommodate extremely bursty traffic overloads, such as 255 ports generating to 1 port, regardless of frame size.

In contrast, the earlier N9K-X9432C-S modules required the addition of 448-usec gaps in test traffic consisting of jumbo frames, to accommodate that module’s smaller buffers. This new round of N-to-N multicast tests includes results both with and without the 448-usec gap to show the effects on latency and jitter.

² For earlier results using the N9K-X9432C-S, see the Network Test report “Cisco Nexus 9508: A New High-Water Mark for 100G Performance.”

³ At test time, hardware version B0 of the N9K-X9432C-S module supported only store-and-forward mode. Cisco has since released hardware version B1 of that module, which supports store-and-forward mode by default, and optionally also allows cut-through mode. Since only hardware version B0 was available for the earlier tests, the Cisco switch used store-and-forward mode in all tests.

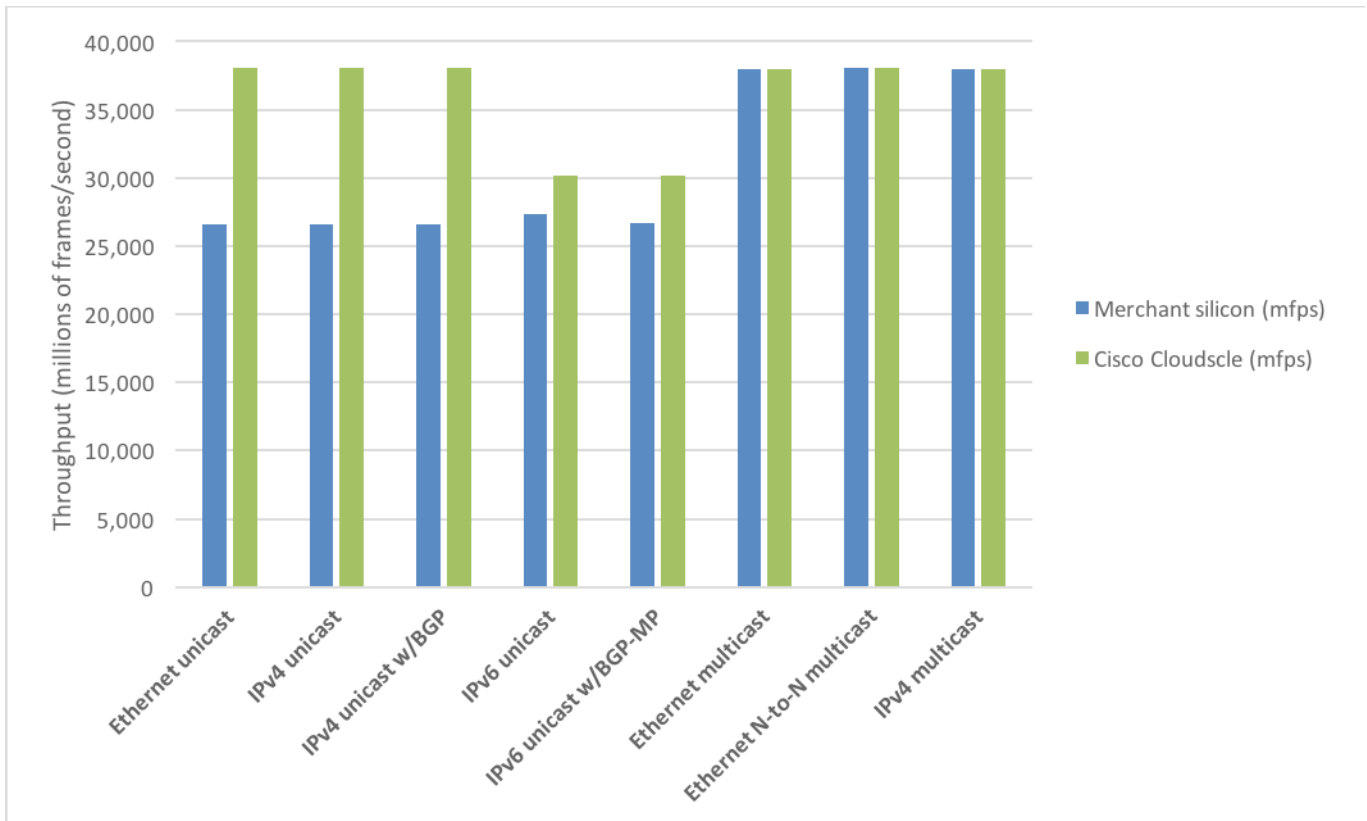


Figure 2: Comparing Cisco Cloudscale and merchant silicon throughput

RFC 2889 Ethernet Unicast Performance

The best way to describe a switch’s forwarding and delay characteristics comes through a test that fully stresses its fabric. [RFC 2889](#) describes such a test, and has long been the industry-standard methodology for Ethernet switch performance testing with unicast traffic.

For this test, engineers configured all 256 100G Ethernet ports of the Cisco switch to be access-mode members of the same VLAN. The Spirent TestCenter traffic generator/analyzer blasted the Cisco Nexus 9508 in a “fully meshed” pattern, meaning test traffic offered to each port was destined to all 255 other ports. The Spirent tool emulated one host attached to each port in the VLAN. In a fully meshed traffic test, the test instrument offers exactly one frame at a time to each destination port; thus, the test traffic pattern does not itself introduce congestion, and any frame loss is a function of fabric capacity.

The test instrument offered traffic at the throughput rate, and also measured latency and jitter at that rate for a variety of frame sizes. Frame sizes range from the Ethernet minimum of 64 bytes to the maximum of 1,518, and beyond to 9,216-byte jumbo frames.

In all tests, the throughput rate was 99.994 percent of line rate. Cisco and Network Test agreed to use 99.994 percent of line rate, which is 60 parts per million (60 ppm) slower than nominal line rate, to avoid clocking differ-



ences between the traffic generator and the switch under test. The IEEE 802.3 Ethernet specification requires interfaces to tolerate clocking differences of up to +/- 100 ppm.

Table 6 presents throughput, latency, and jitter results from the Ethernet unicast tests. Note that throughput for the Cisco Nexus 9508 was essentially line rate in every test case.

For this and all switch tests described here, engineers configured the Cisco device to operate in “cut-through” mode, where the switch begins forwarding each incoming frame before it has fully cached the frame. As a general rule, cut-through mode offers the lowest possible latency and jitter, and also usually delivers the same predictable average latency across all frame sizes.

This is because the [RFC 1242](#), the terminology companion document to RFC 2544, requires cut-through devices to be tested using first-in, first-out (FIFO) latency measurements. With cut-through mode, the first byte of each frame appears on an output interface at the same time, regardless of frame size. In contrast, RFC 1242 requires latency measurement of store-and-forward devices to use a last-in, first-out (LIFO) method. Here, average latency will increase roughly proportional to frame size, since the device must cache the entire frame before forwarding it.

Figures 3 and 4 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.994 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,092,952,698	19.504	99.994%	0.970	2.850	4.290	0.009	1.140
128	21,620,324,502	22.139	99.994%	0.960	2.878	3.710	0.007	0.430
256	11,593,507,325	23.744	99.994%	0.970	2.929	3.930	0.007	0.500
512	6,014,676,717	24.636	99.994%	0.990	3.025	5.090	0.007	1.180
1,024	3,064,950,207	25.108	99.994%	1.010	3.041	4.550	0.005	1.150
1,280	2,461,390,780	25.205	99.994%	1.000	3.042	5.310	0.005	1.780
1,518	2,080,499,362	25.266	99.994%	1.010	3.028	4.800	0.006	1.410
9,216	346,449,552	25.543	99.994%	1.020	2.993	5.280	0.008	1.970

Table 6: Ethernet unicast performance results

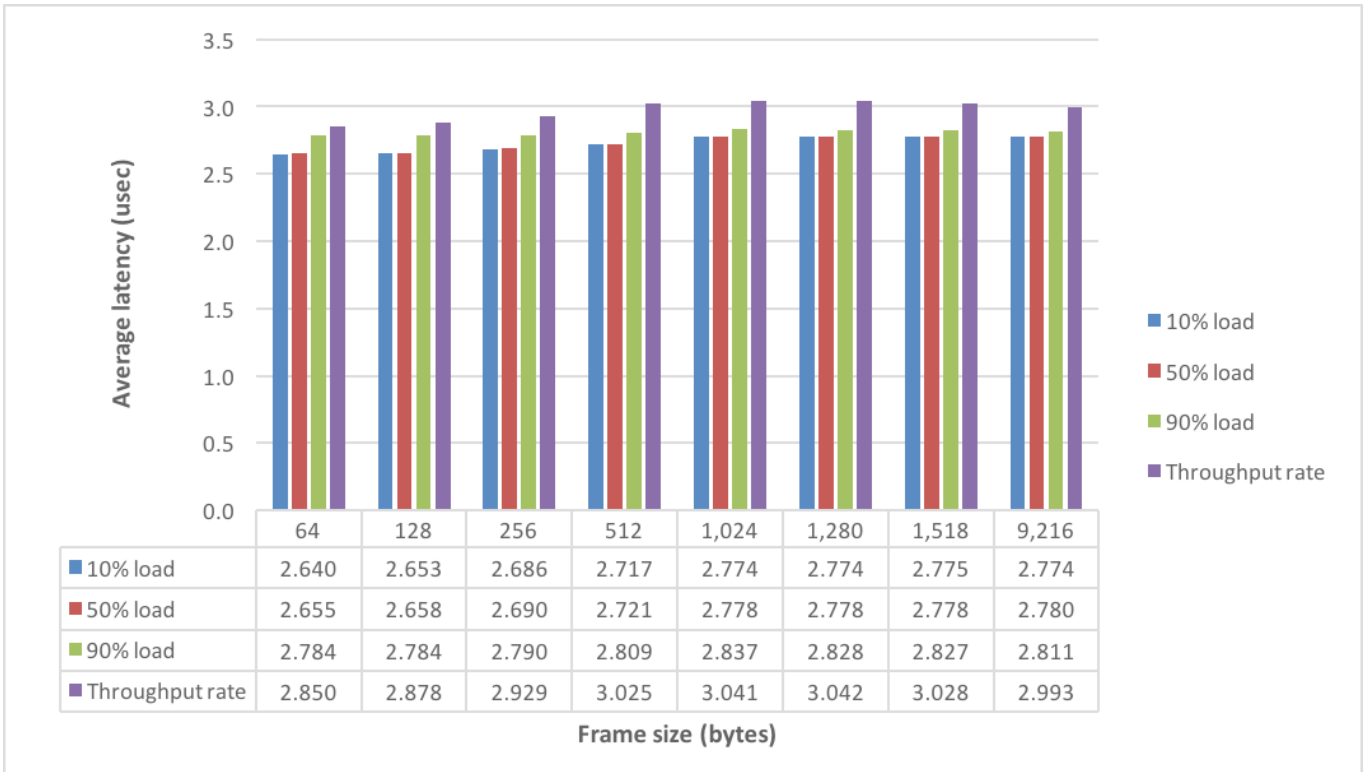


Figure 3: Ethernet unicast average latency vs. load

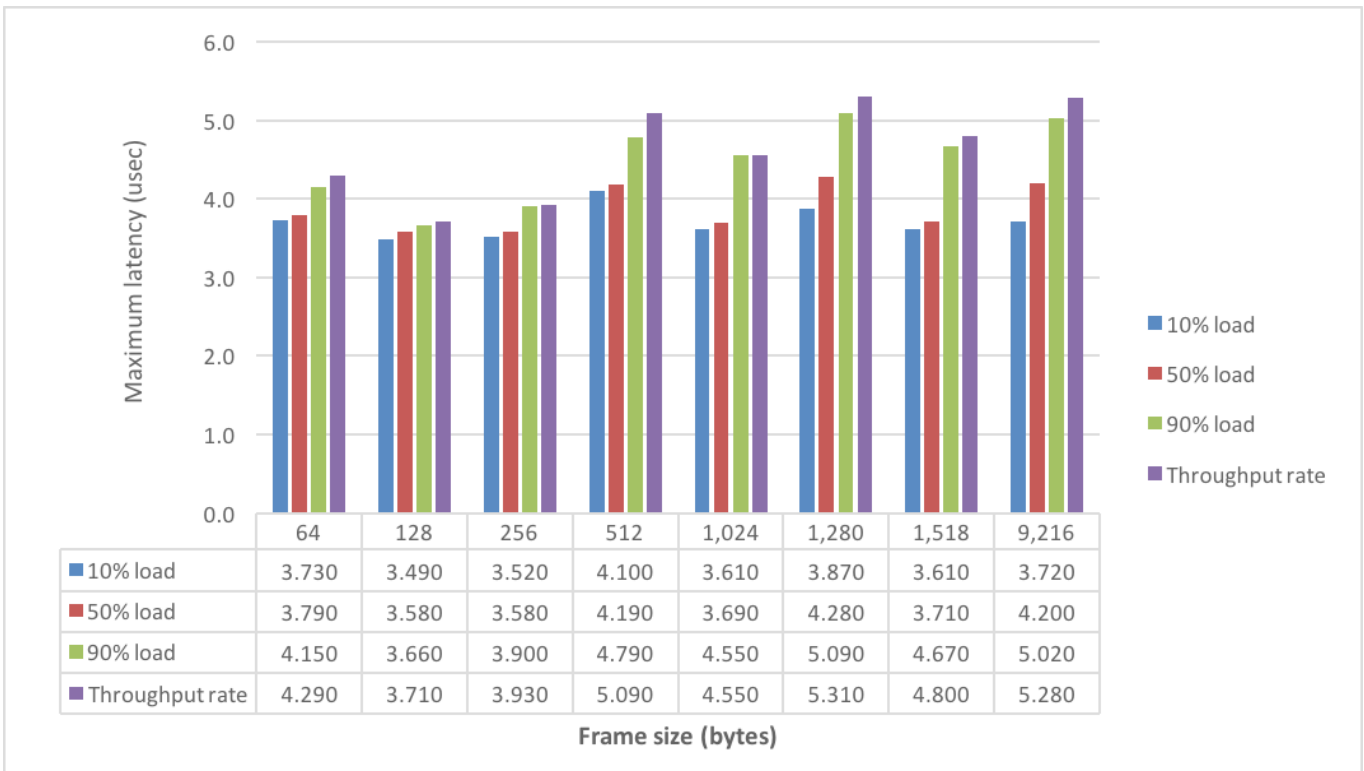


Figure 4: Ethernet unicast maximum latency vs. load



RFC 2544 IPv4 Unicast Performance

The Cisco Nexus 9508 acted as a router in the IPv4 performance tests, with each interface on a different IP subnet. As in the Ethernet tests, engineers used fully meshed traffic among all ports to measure throughput, latency, and jitter. The Spirent test instrument again emulated one host per subnet.

As noted in the “Comparing Cisco and Merchant Silicon ASICs” section, test traffic carried UDP as well as IP headers.

Throughput again was equivalent to 99.994 percent of line rate in all tests, regardless of frame size. Table 7 presents throughput, latency, and jitter results from the IPv4 unicast tests.

Figures 5 and 6 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.994 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,092,952,660	19.504	99.994%	0.950	2.920	4.300	0.009	1.140
128	21,620,324,473	22.139	99.994%	0.960	2.938	3.820	0.008	0.490
256	11,593,507,314	23.744	99.994%	0.970	2.992	3.920	0.007	0.550
512	6,014,676,714	24.636	99.994%	0.970	3.087	4.830	0.007	1.070
1,024	3,064,950,204	25.108	99.994%	1.000	3.111	4.520	0.004	1.180
1,280	2,461,390,779	25.205	99.994%	0.990	3.115	5.620	0.005	1.880
1,518	2,080,499,360	25.266	99.994%	0.990	3.100	4.780	0.006	1.450
9,216	346,449,552	25.543	99.994%	0.990	3.058	5.750	0.007	2.010

Table 7: IPv4 unicast performance results

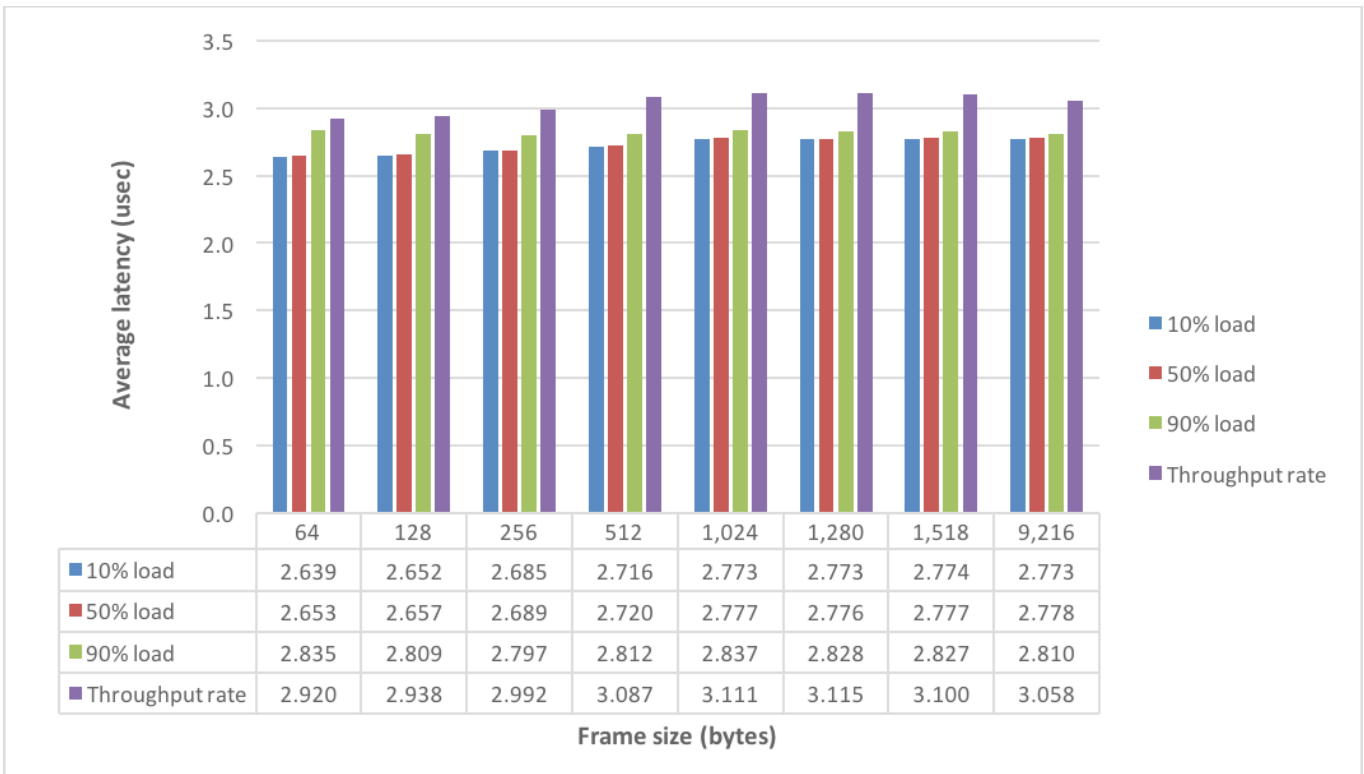


Figure 5: IPv4 unicast average latency vs. load

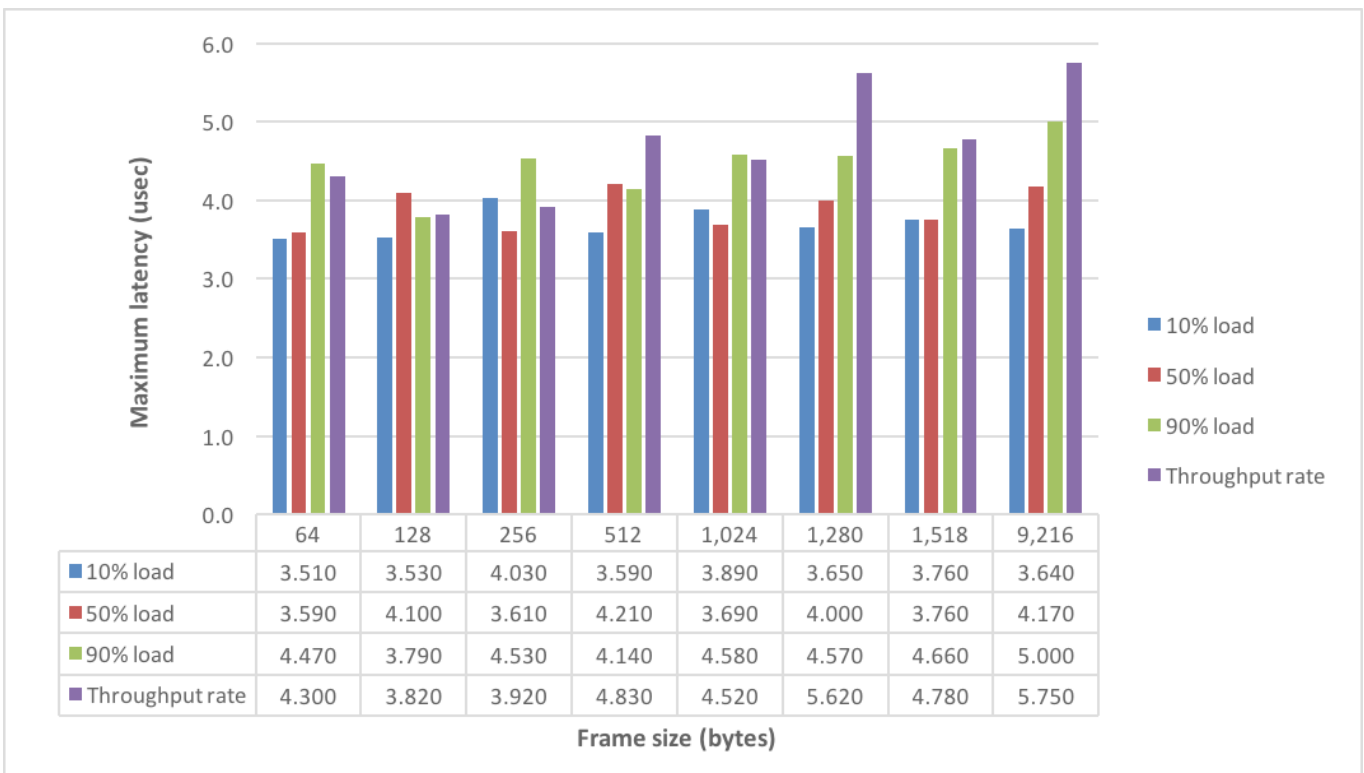


Figure 6: IPv4 unicast maximum latency vs. load



RFC 2544 IPv4 Unicast Performance With BGP Routing

While IPv4 tests moved traffic between local configured subnets, they used only direct, locally configured routes instead of dynamic routes. Engineers also assessed IPv4 performance using Border Gateway Protocol (BGP), in this case moving traffic among 2,048 unique networks learned via BGP.

Test engineers configured Spirent TestCenter to emulate 256 BGP routing peers, each using a unique Autonomous System Number (ASN). Each Spirent BGP router brought up a peering session with the Cisco Nexus 9508, then advertised a total of 2,048 unique routes. The Spirent test tool then offered fully meshed traffic between all networks learned using BGP.

Note that the choice of 2,048 routes is due to a limit in the number of trackable receive streams supported by the Spirent dX2 test modules. Cisco says the Nexus 9508 equipped with N9K-X9732C-EX modules supports up to 736,000 longest-prefix match (LPM) routes and up to 736,000 host entries, but Network Test did not verify this. Other Spirent test modules also support higher trackable stream counts. With the dX2 module, a higher route count also would have been possible using fewer than 256 ports.

Throughput again was equivalent to 99.994 percent of line rate in all tests, regardless of frame size, despite the higher amount of overall traffic due to BGP control-plane traffic.

Table 8 presents throughput, latency, and jitter results for all frame sizes.

Figures 7 and 8 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.994 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,092,948,084	19.504	99.994%	0.960	3.033	4.960	0.006	1.650
128	21,620,321,875	22.139	99.994%	0.960	3.024	4.720	0.007	1.280
256	11,593,507,286	23.744	99.994%	0.970	3.011	5.140	0.006	1.800
512	6,014,676,715	24.636	99.994%	0.970	3.028	6.180	0.006	2.080
1,024	3,064,950,205	25.108	99.994%	1.000	3.091	5.550	0.007	2.130
1,280	2,461,390,782	25.205	99.994%	0.990	3.077	5.980	0.006	2.070
1,518	2,080,499,361	25.266	99.994%	0.990	3.087	5.520	0.007	2.100
9,216	346,449,551	25.543	99.994%	1.000	3.050	5.260	0.007	1.980

Table 8: IPv4/BGP unicast performance results

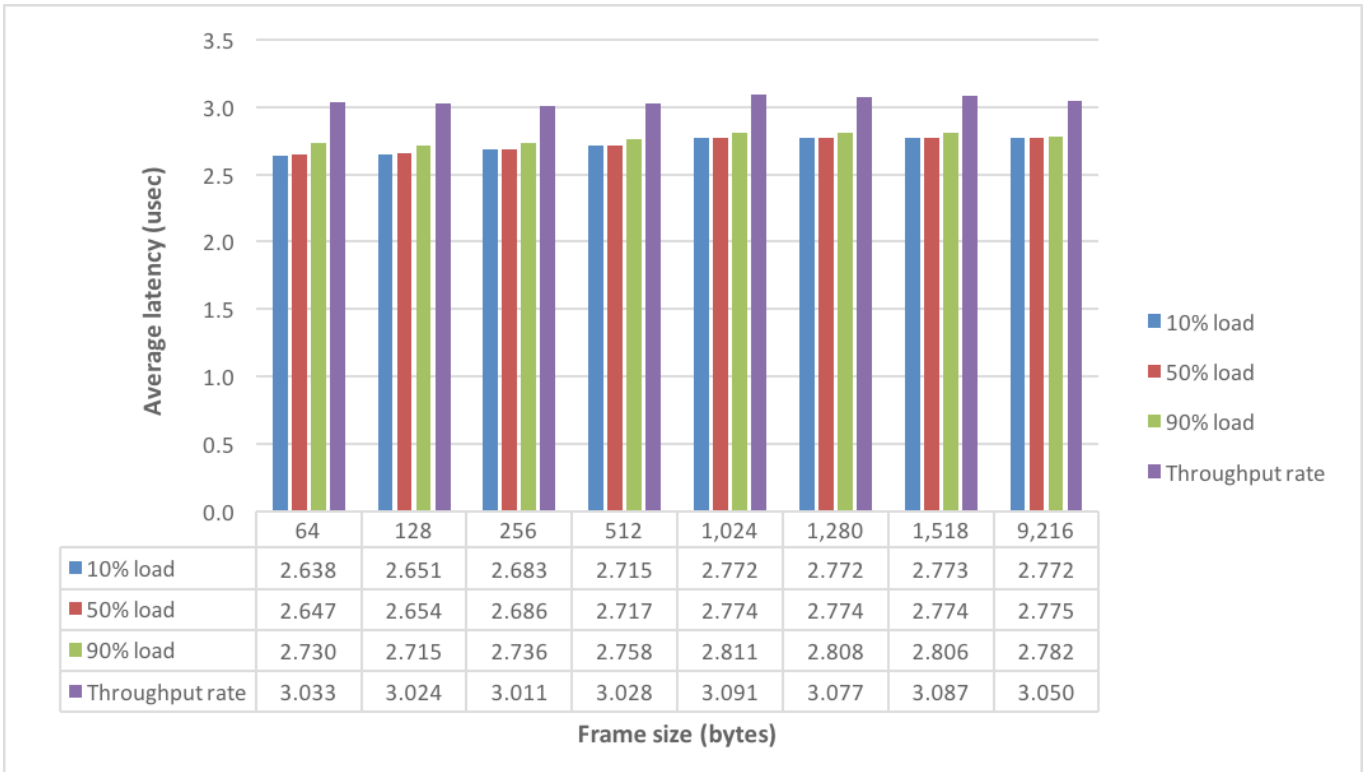


Figure 7: IPv4/BGP unicast average latency vs. load

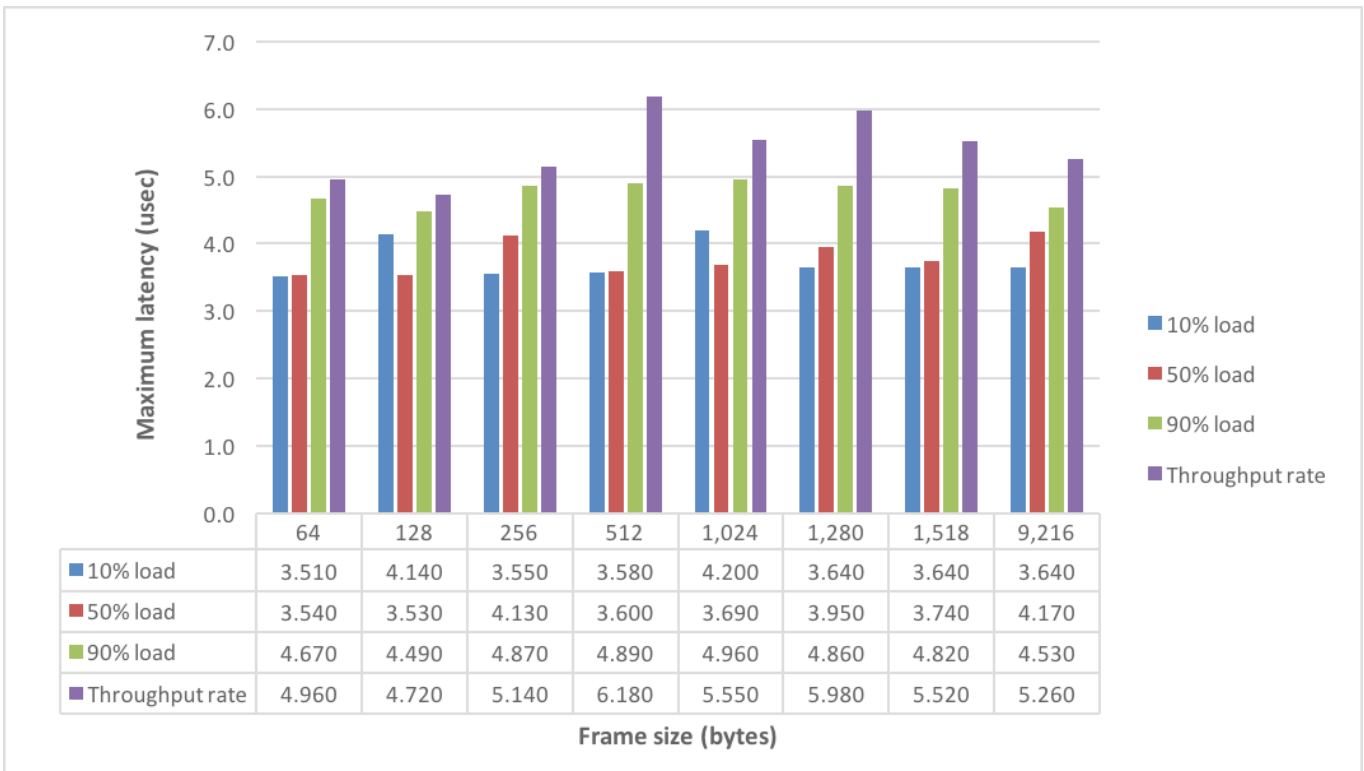


Figure 8: IPv4/BGP unicast maximum latency vs. load



RFC 5180 IPv6 Unicast Performance

As in previous routing tests, IPv6 performance measurements required the Cisco Nexus 9508 to route traffic among unique subnets on each of 256 interfaces. Test traffic again was fully meshed, meaning traffic offered to each port went to all other ports. The Spirent test instrument emulated one IPv6 host per subnet.

For all IPv6 tests, test engineers configured a minimum frame size of 86 bytes rather than 64 bytes to accommodate the 20-byte “signature field” added by the Spirent test instrument, plus an 8-byte UDP header (see the Test Methodology section for more details on 86-byte frames). Test traffic again used UDP headers for comparison with earlier test results, as discussed in the “Comparing Cisco Cloudscale Technology and Merchant Silicon ASICs” section.

Throughput again was equivalent to 99.994 percent of line rate in most tests. Table 9 presents throughput, latency, and jitter results for all frame sizes.

Figures 9 and 10 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.994 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
86	30,186,868,754	20.769	99.994%	2.210	2.964	3.730	0.009	0.520
128	21,620,324,920	22.139	99.994%	2.230	2.970	4.370	0.008	1.090
256	11,593,507,575	23.744	99.994%	2.260	3.003	3.850	0.007	0.550
512	6,014,676,849	24.636	99.994%	2.290	3.093	4.130	0.007	0.640
1,024	3,064,950,275	25.108	99.994%	2.360	3.118	4.970	0.004	1.170
1,280	2,461,390,839	25.205	99.994%	2.350	3.098	4.280	0.005	0.630
1,518	2,080,499,408	25.266	99.994%	2.350	3.104	4.830	0.006	1.370
9,216	346,449,560	25.543	99.994%	2.370	3.065	3.980	0.008	0.480

Table 9: IPv6 unicast performance results

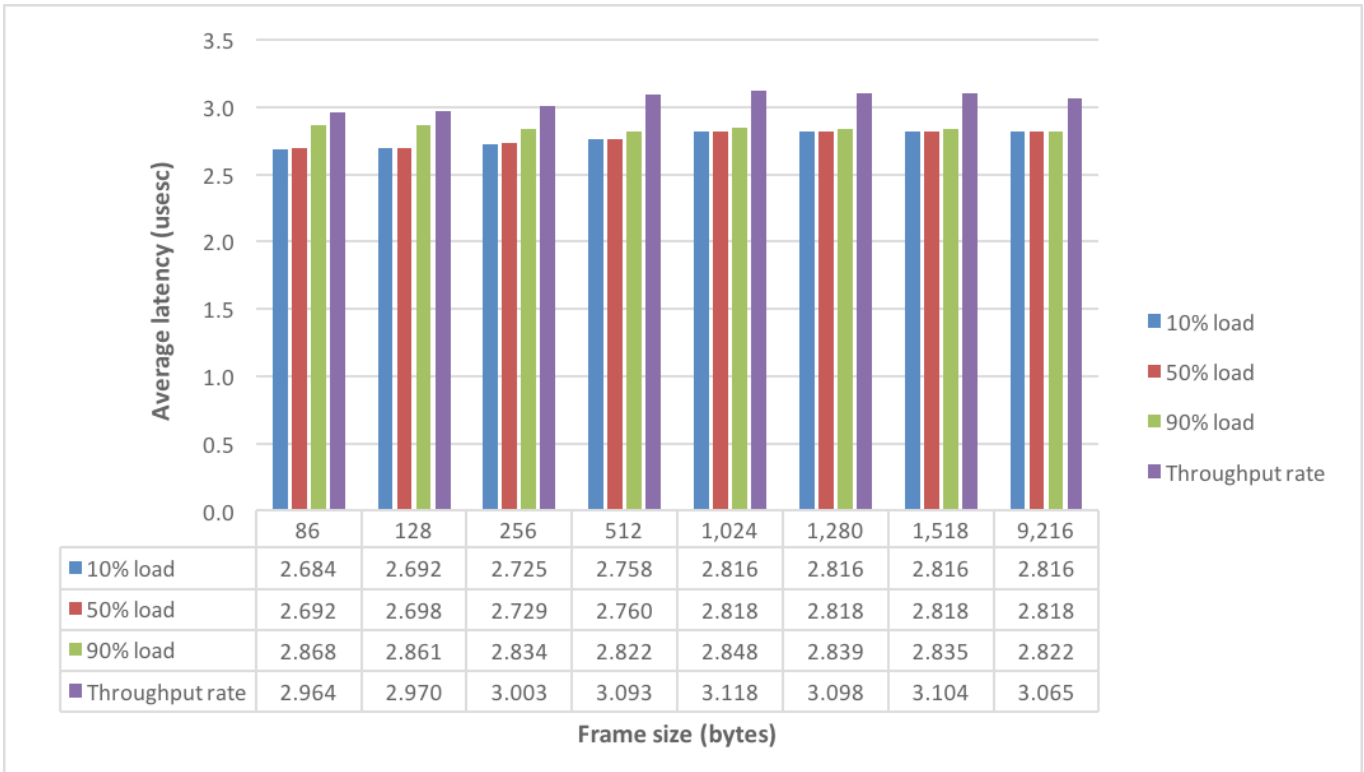


Figure 9: IPv6 unicast average latency vs. load

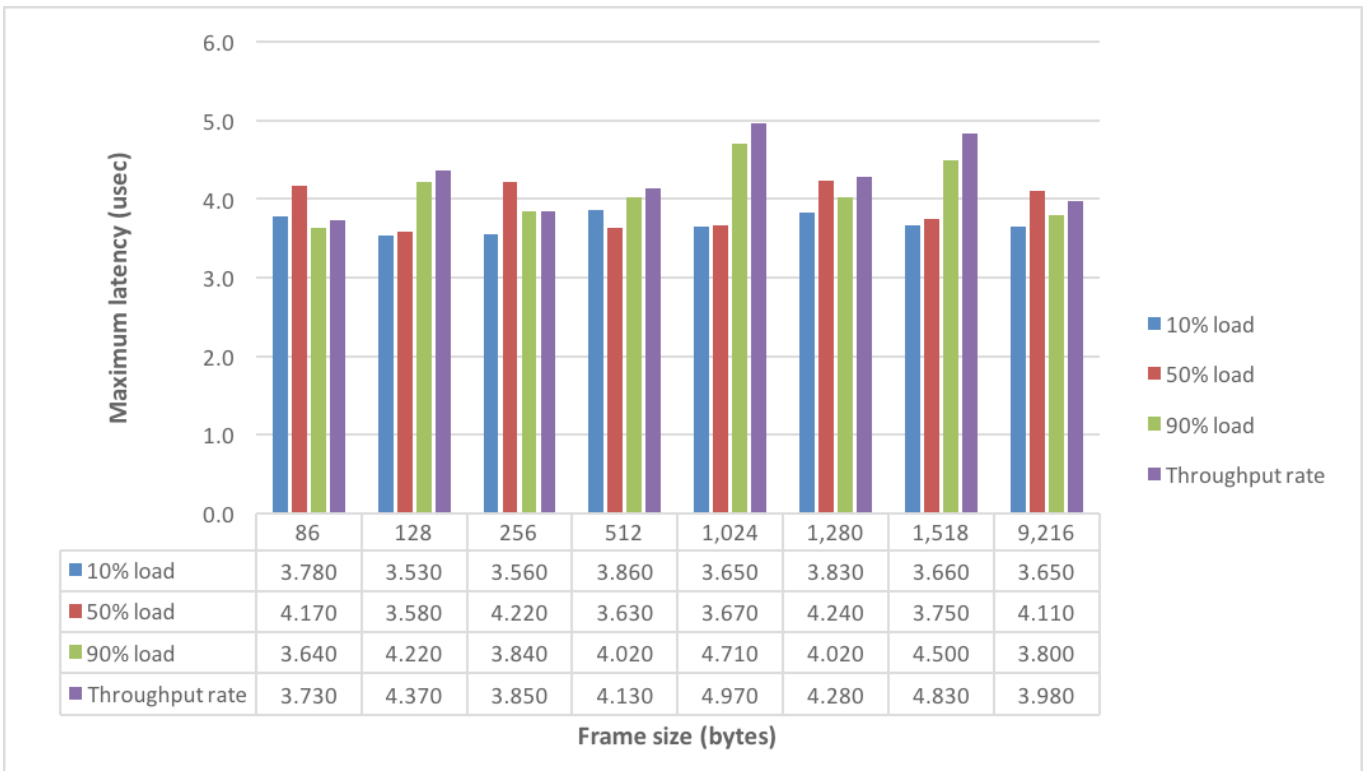


Figure 10: IPv6 unicast maximum latency vs. load



RFC 5180 IPv6 Unicast Performance With BGP-MP Routing

Just as IPv4 tests involved cases without and then with BGP routing, engineers also tested IPv6 with dynamic routing enabled. This tests adds BGP-Multiprotocol (BGP-MP) routing to the previous IPv6 test. Here, the Cisco Nexus 9508 routed traffic among 2,048 unique IPv6 networks learned via BGP-MP.

[RFC 4760](#) describes BGP-MP, a set of multiprotocol extensions to BGP for carrying topology information about different network-layer protocols, including IPv6. The Spirent test instrument brought up BGP-MP peering sessions on each port, advertised 2,048 unique IPv6 routes, and then offered fully meshed traffic destined to all routes.

Note that the choice of 2,048 routes is due to a limit in the number of trackable receive streams supported by the Spirent dX2 test modules. Cisco says the Nexus 9508 equipped with N9K-X9732C-EX modules supports up to 234,000 longest-prefix match (LPM) routes and 34,000 host entries according to the Cisco data sheet, but Network Test did not verify this. Other Spirent test modules also support higher trackable stream counts. A higher route count also would have been possible using fewer than 256 ports.

The minimum frame size in these tests was 86 bytes for the reasons discussed in the “IPv6 Unicast Performance” section.

Throughput again was equivalent to 99.994 percent of line rate for all tests, despite the higher amount of overall traffic due to BGP-MP control-plane traffic.

Table 10 presents throughput, latency, and jitter results for all frame sizes. Figures 11 and 12 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.994 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
86	30,186,867,941	20.769	99.994%	2.220	3.034	4.690	0.007	0.890
128	21,620,324,308	22.139	99.994%	2.230	3.057	4.140	0.006	0.700
256	11,593,507,262	23.744	99.994%	2.260	3.022	4.120	0.006	0.660
512	6,014,676,692	24.636	99.994%	2.290	3.050	4.800	0.006	1.350
1,024	3,064,950,198	25.108	99.994%	2.370	3.093	4.300	0.006	0.590
1,280	2,461,390,775	25.205	99.994%	2.360	3.102	5.180	0.006	1.700
1,518	2,080,499,358	25.266	99.994%	2.360	3.087	4.130	0.007	0.580
9,216	346,449,551	25.543	99.994%	2.350	3.072	4.580	0.007	0.650

Table 10: IPv6/BGP-MP unicast performance results

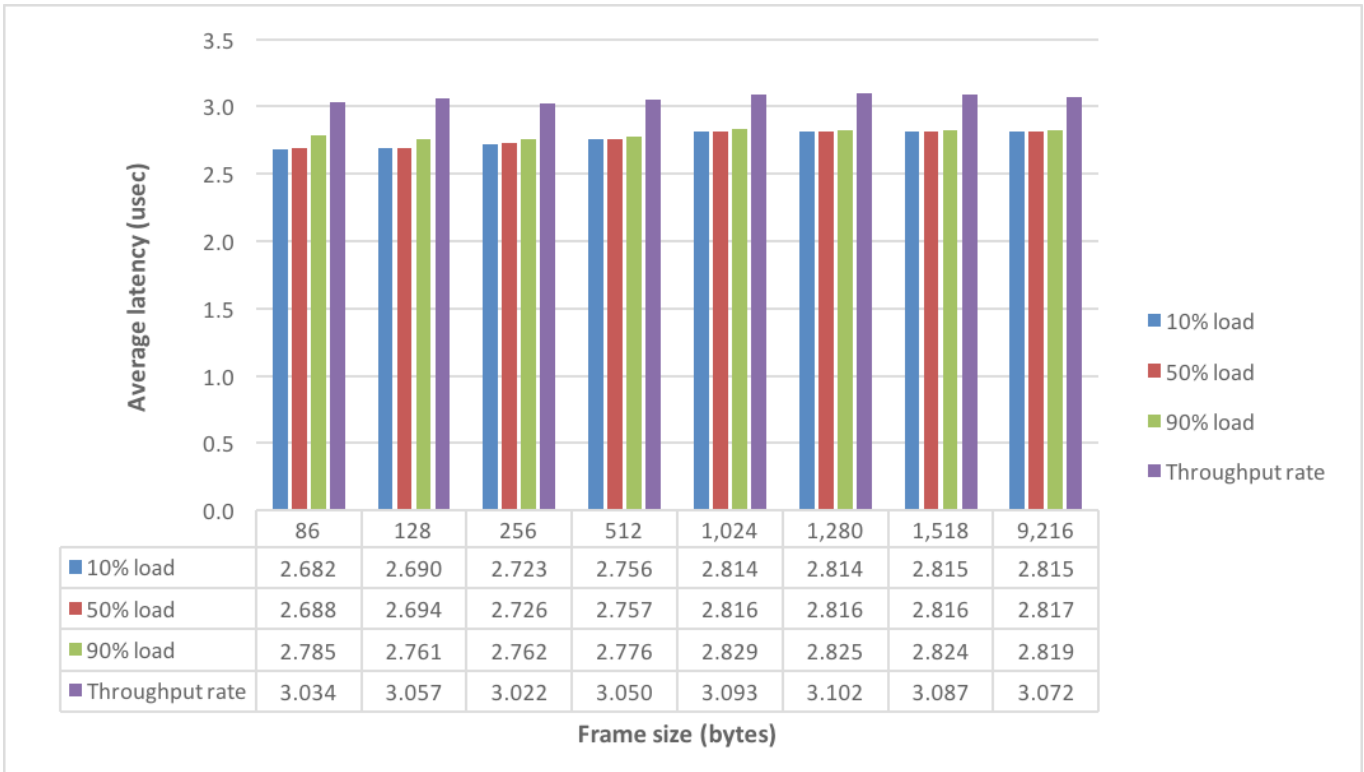


Figure 11: IPv6/BGP-MP unicast average latency vs. load

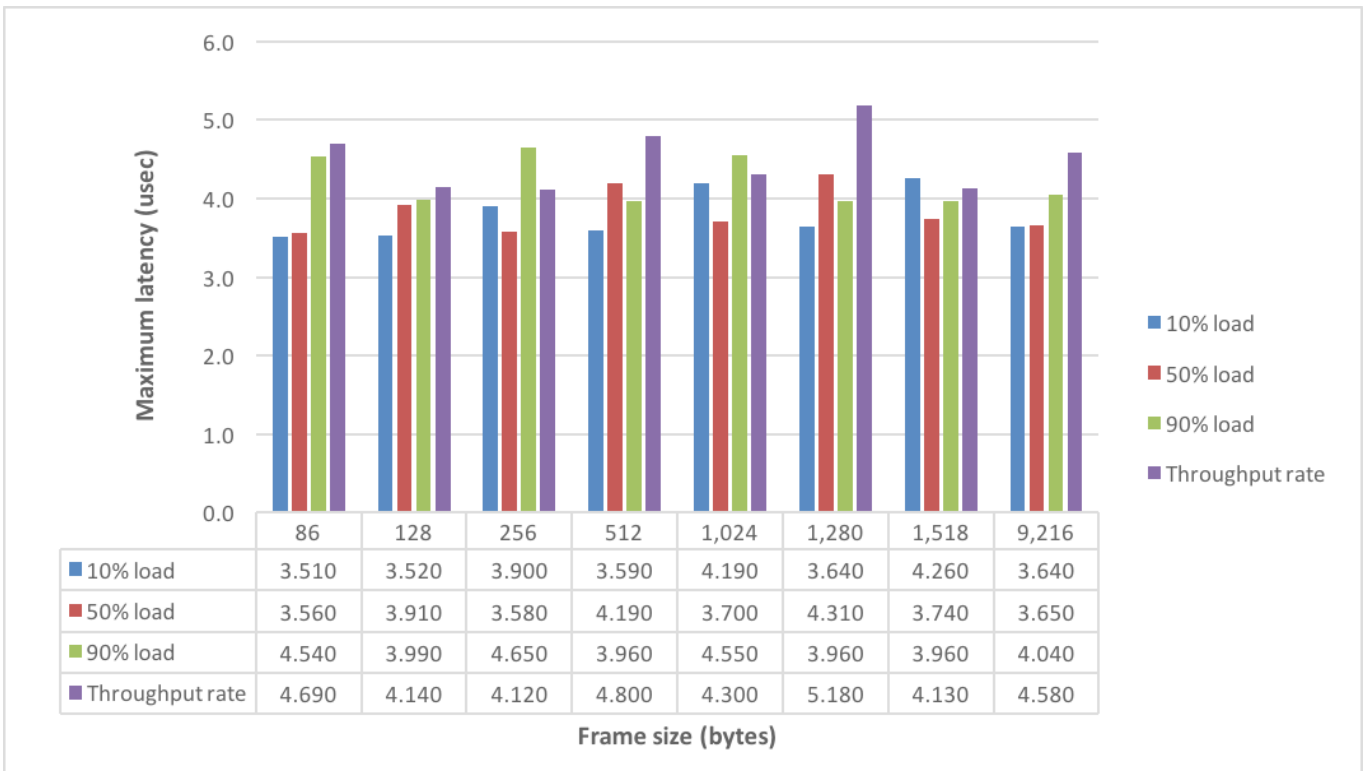


Figure 12: IPv6/BGP-MP unicast maximum latency vs. load



RFC 3918 Ethernet 1-to-N Multicast Performance

With IP multicast, a single incoming frame may be replicated to all other switch ports (255 in this case). Multicast tests also involve both control and data planes, the former for group subscription and the latter for traffic forwarding. Thus, multicast tests of the Cisco Nexus 9508 were highly stressful, and scaled up the switch to its maximum performance levels.

On the control plane, the Spirent test instrument emulated hosts on 255 subscriber ports, each joining the same 4,095 IP multicast groups, thus requiring the Cisco device to maintain forwarding information for more than 1 million unique entities. This is calculated as 4,095 multicast routes (mroutes) times an outgoing interface list (OIL) of 255 ports for a total mroute*OIL count of 1,044,225.

A multicast test with more than 1 million mroutes*OIL entries is significant; it represents the largest Ethernet multicast table ever constructed using Cisco Nexus 9000 technology.

On the data plane, the Spirent instrument offered traffic in a way that required massive replication on the part of the Cisco switch. The Spirent generator offered traffic destined to all 4,095 IP multicast group addresses on all 255 receiver interfaces. The use of 4,095 multicast groups was due to a stream count limit in the Spirent dX2 test module, and is not a limit of the Cisco Nexus 9508.

In this Ethernet test, all interfaces on the Cisco switch were members of a single VLAN. The switch used IGMP version 3 (IGMPv3) to build IGMP snooping tables.

The Cisco Nexus 9508 delivered all multicast traffic with zero frame loss with throughput equivalent to 99.994 percent of line rate in all test cases. Table 11 presents throughput, latency, and jitter results for all frame sizes.

Figures 13 and 14 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.994 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	37,944,152,063	19.427	99.994%	1.240	3.002	3.630	0.007	0.700
128	21,535,870,094	22.053	99.994%	1.250	3.014	3.370	0.005	0.070
256	11,548,220,191	23.651	99.994%	1.250	3.043	3.400	0.004	0.060
512	5,991,181,903	24.540	99.994%	1.260	3.074	4.050	0.004	0.650
1,024	3,052,977,752	25.010	99.994%	1.280	3.130	3.460	0.004	0.040
1,280	2,451,775,979	25.106	99.994%	1.280	3.136	3.470	0.004	0.040
1,518	2,072,372,415	25.167	99.994%	1.280	3.141	4.160	0.004	0.620
9,216	345,096,232	25.443	99.994%	1.280	3.370	3.710	0.004	0.050

Table 11: Ethernet 1-to-N multicast performance results

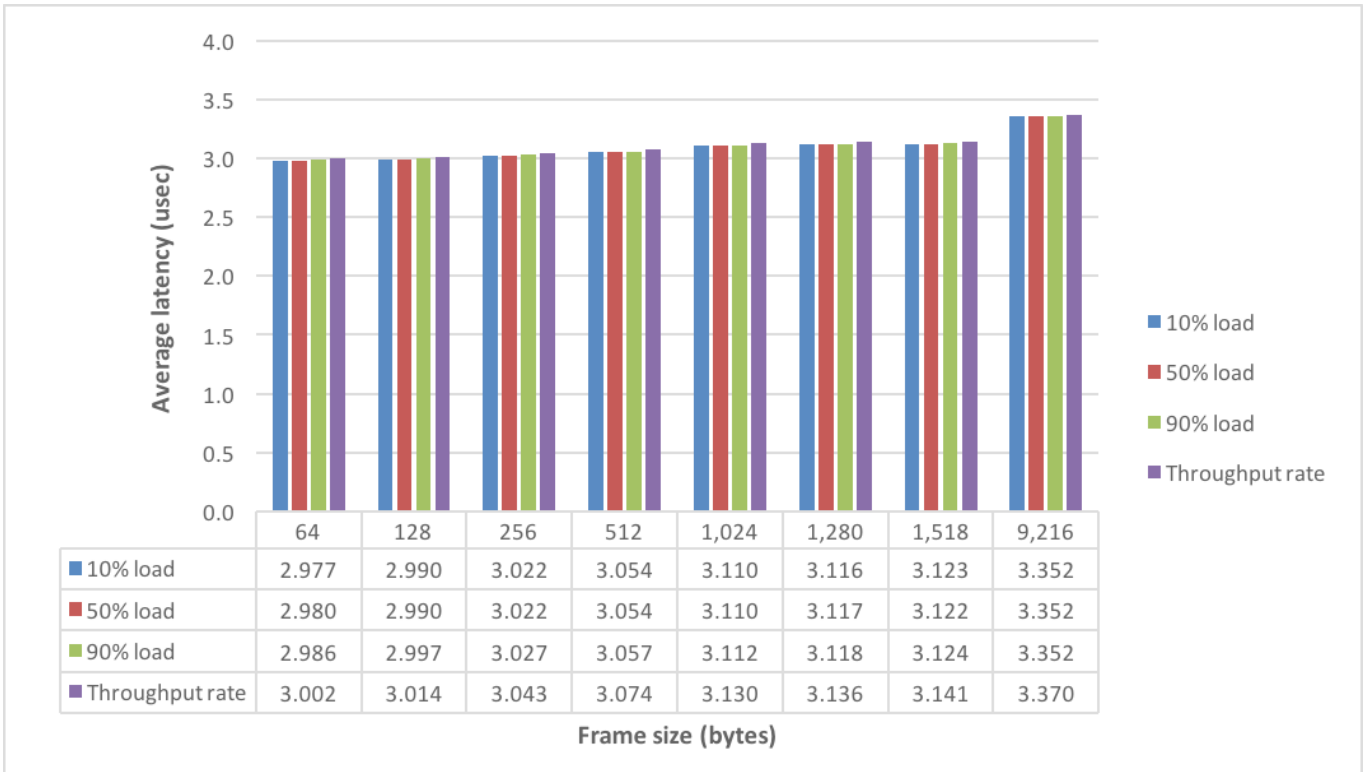


Figure 13: Ethernet 1-to-N multicast average latency vs. load

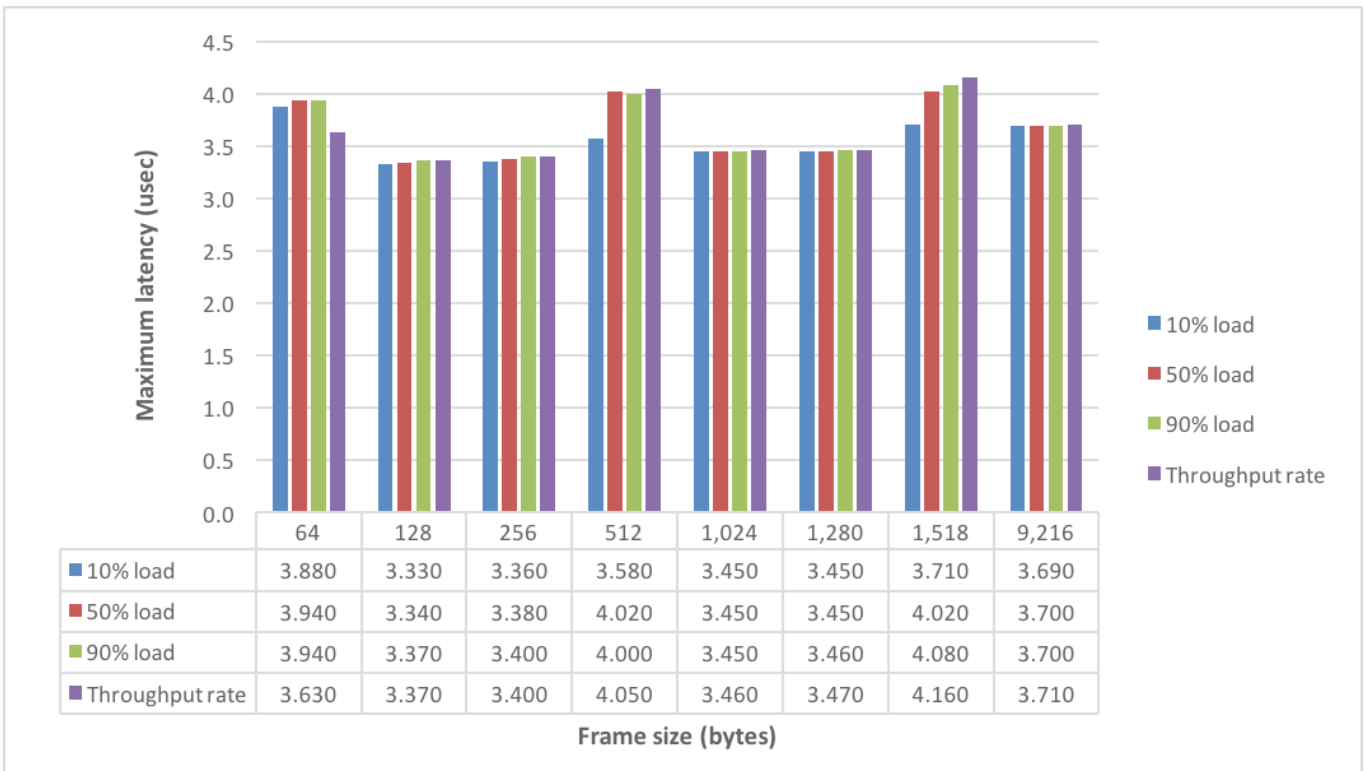


Figure 14: Ethernet 1-to-N multicast maximum latency vs. load



Ethernet N-to-N Multicast Performance

Because industry-standard [RFC 3918](#) multicast tests require only one multicast transmitter, they do not meaningfully predict performance of applications that involve many concurrent transmitters and receivers. A high-profile example is trading-floor applications, where financial institutions subscribe to multiple stock-quote feeds, each transmitting to a different multicast group address, and all sending traffic to multiple traders.

To model this kind of N-to-N multicast connectivity, Cisco and Network Test devised a benchmark in which all 256 100G Ethernet ports are both multicast transmitters and receivers.

In this scenario, engineers configured each port of the Spirent TestCenter instrument to offer traffic to one multicast group address, with the remaining 255 ports subscribed to that multicast group. With all ports configured this way, each Spirent port offered traffic to one multicast group address, and expected to receive traffic from the 255 other multicast group addresses on 255 other ports.

By definition, this was an overload test, since each output port received 255 frames at the same instant. In switch modules with relatively small buffers, the output ports may not have sufficient capacity to handle a 255:1 overload. Indeed, in previous tests involving the N9K-X9432C-S module, the switch would drop jumbo frames, regardless of load, due to the overload pattern. Since jumbo frames are seldom used in multicast applications, this was an edge case.

Fortunately, the new N9K-X9732C-EX modules for the Cisco Nexus 9508 have higher buffer capacity than previous products, and can handle even this edge case without frame loss. The tradeoff: Latency will naturally increase because of the larger buffers.

As noted, in previous N-to-N tests with N9K-X9432C-S modules, the switch dropped jumbo frames at any load, making it impossible to measure throughput and latency. Instead of tuning switch buffers to non-default sizes, engineers instead configured Spirent TestCenter to use a “staggered start” pattern. Normally, the Spirent test instrument begins to transmit each frame on all ports at the same instant; in this case, engineers configured a staggered-start interval of 448 usec between ports to avoid an overload in the jumbo-frame test case. The Spirent instrument sets asynchronous traffic patterns in 64-usec intervals. Through trial and error, test engineers determined that 7 intervals, or 448 usec, was the minimum needed to conduct the earlier test with zero loss.

This staggered-start pattern isn’t necessary with N9K-X9732C-EX modules, since they can forward N-to-N traffic for any frame size with zero loss. However, since previous tests *did* use a staggered start, test engineers this time measured N-to-N performance twice, both with and without the staggered-start pattern.

As in the 1-to-N multicast tests, the Cisco Nexus 9508 delivered N-to-N multicast traffic with zero frame loss in all test cases. Table 12 presents throughput, latency, and jitter results for all frame sizes with a synchronous start (with all traffic beginning simultaneously). Table 13 presents performance results with a staggered start (with traffic on different ports beginning asynchronously at 448-usec intervals).

Note that latency and jitter with staggered-start, asynchronous traffic is far lower than with synchronous traffic.

Figures 15 and 16 compare average and maximum delay measurements with a synchronous start, using offered loads of 10, 50, 90, and the throughput rate. Figures 17 and 18 compare average and maximum delay measurements with a staggered start, again using offered loads of 10, 50, 90, and the throughput rate.



Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,092,629,392	19.503	99.993%	1.030	3.498	4.690	0.115	1.570
128	21,620,141,019	22.139	99.993%	1.030	4.100	6.460	0.104	2.860
256	11,593,408,965	23.743	99.993%	1.060	5.360	8.420	0.160	3.130
512	6,014,625,719	24.636	99.993%	1.090	7.848	13.520	0.282	6.320
1,024	3,064,924,232	25.108	99.993%	1.150	12.821	23.690	0.566	17.800
1,280	2,461,369,927	25.204	99.993%	1.160	15.285	28.750	0.662	23.820
1,518	2,080,481,739	25.265	99.993%	1.190	17.575	33.470	1.215	29.790
9,216	346,446,621	25.543	99.993%	1.050	95.149	189.110	13.840	181.060

Table 12: Ethernet N-to-N multicast performance results with synchronous traffic

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,020,239,928	19.466	99.803%	1.010	2.879	3.740	0.008	0.100
128	21,579,055,556	22.097	99.803%	1.010	2.906	3.770	0.007	0.110
256	11,571,378,118	23.698	99.803%	1.030	3.026	4.360	0.014	0.690
512	6,003,196,685	24.589	99.803%	1.030	3.059	3.940	0.009	0.270
1,024	3,059,100,758	25.060	99.803%	1.060	3.386	4.400	0.007	0.280
1,280	2,456,693,437	25.157	99.803%	1.050	3.232	4.450	0.007	0.720
1,518	2,076,529,080	25.217	99.803%	1.090	3.623	4.910	0.036	1.490
9,216	345,789,310	25.494	99.803%	1.050	4.381	6.990	0.008	1.300

Table 13: Ethernet N-to-N multicast performance results with asynchronous traffic

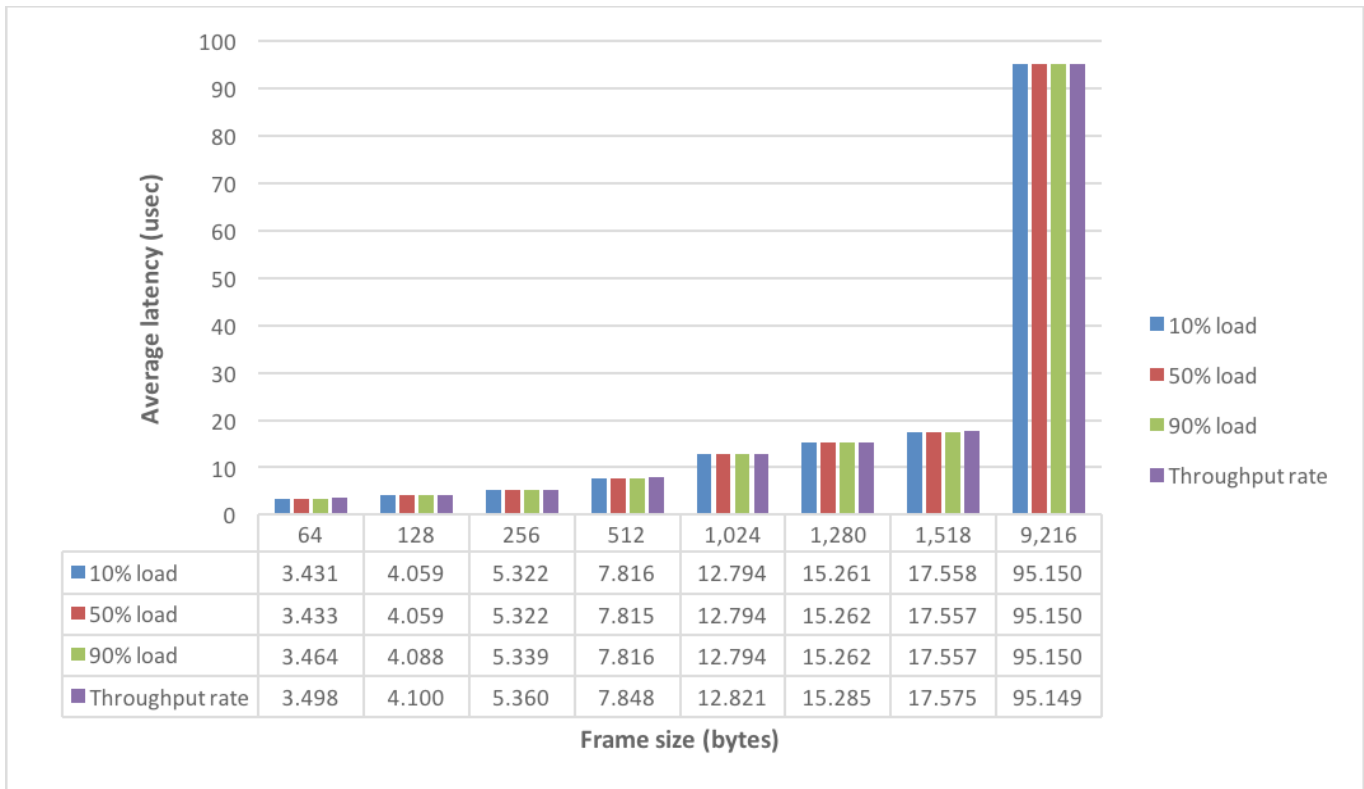


Figure 15: Ethernet N-to-N multicast average latency vs. load with synchronous traffic

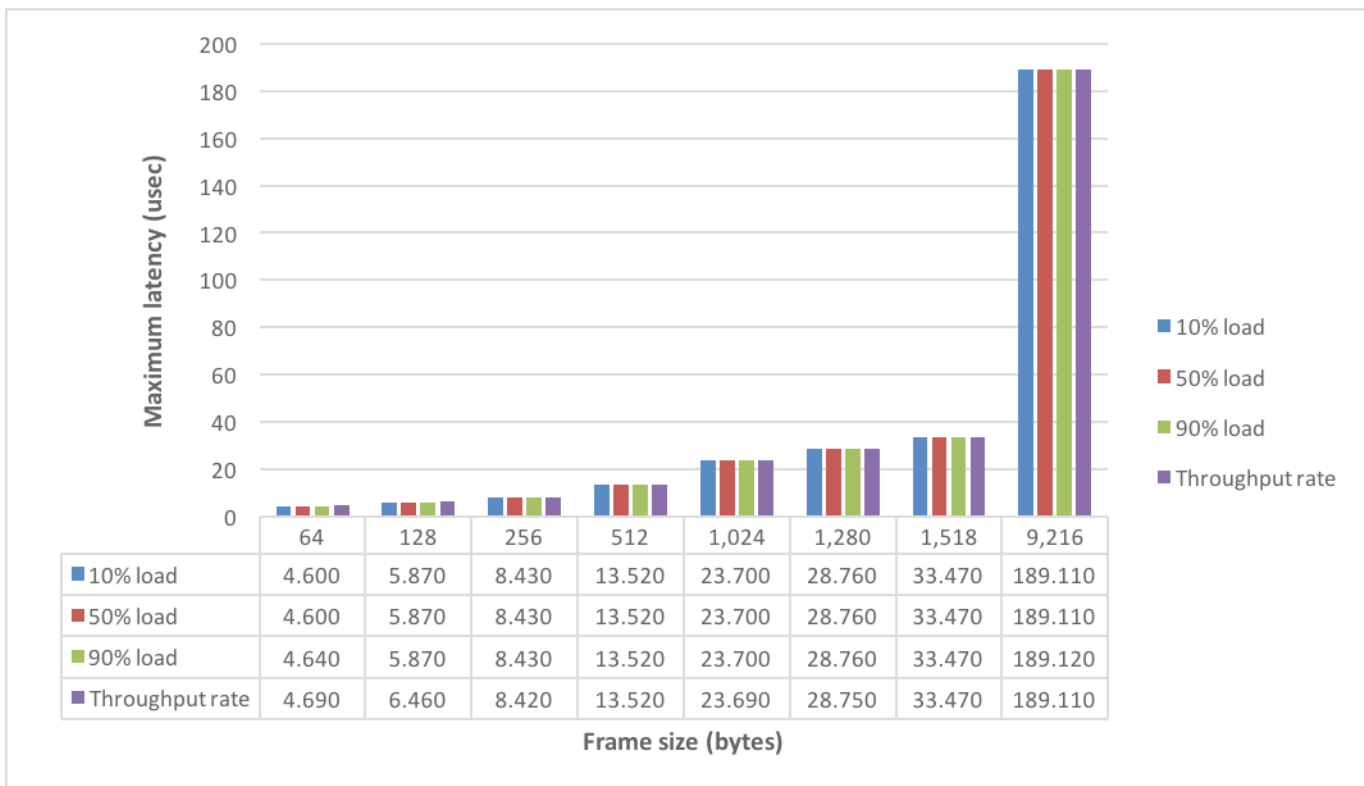


Figure 16: Ethernet N-to-N multicast maximum latency vs. load with synchronous traffic

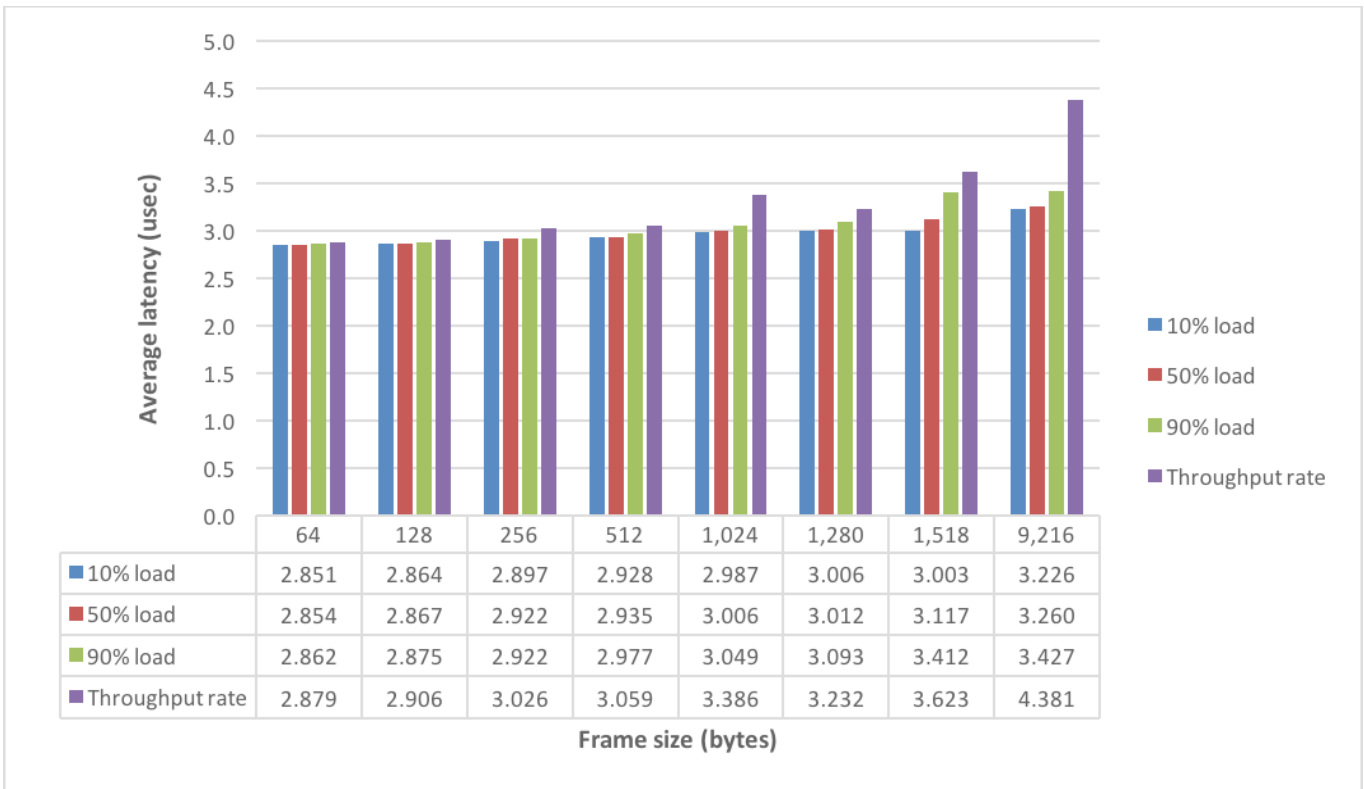


Figure 17: Ethernet N-to-N multicast average latency vs. load with asynchronous traffic

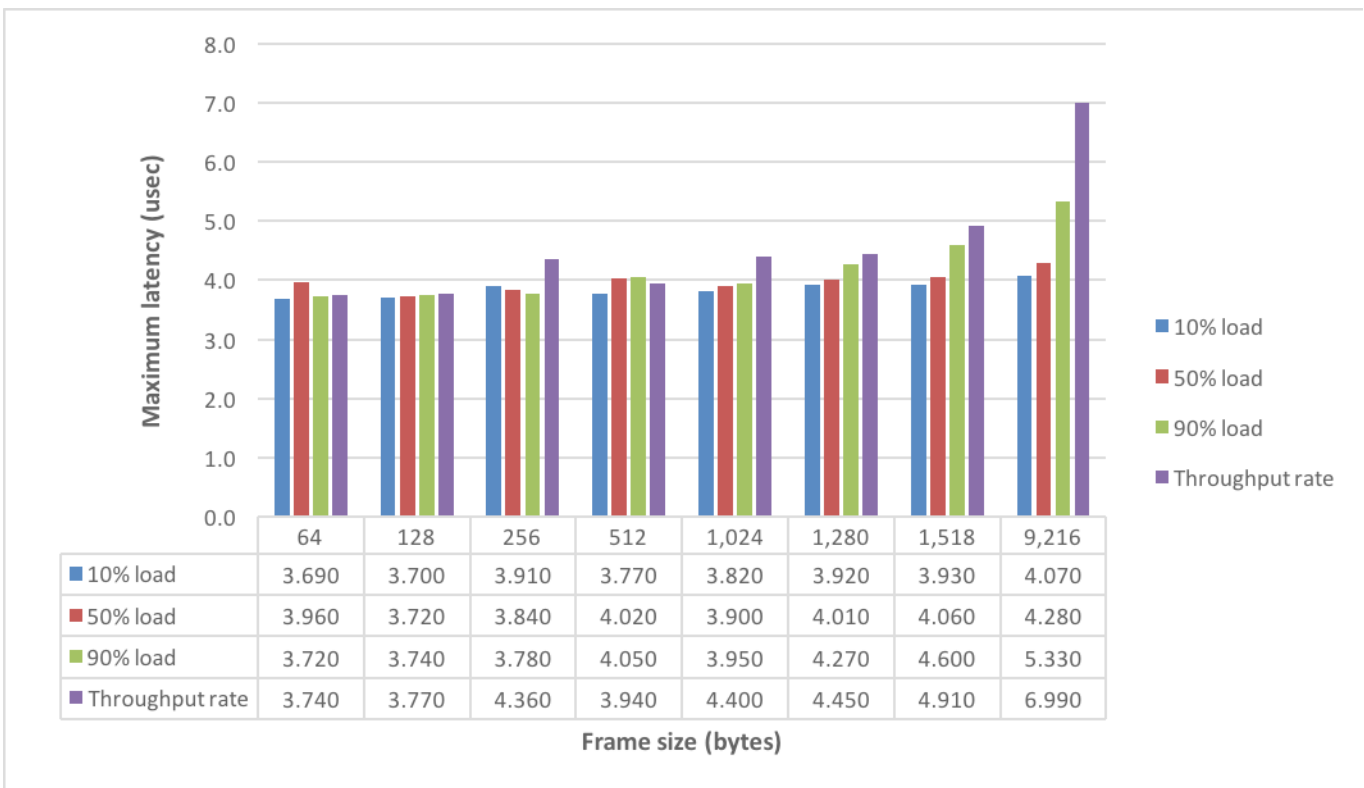


Figure 18: Ethernet N-to-N multicast maximum latency vs. load with asynchronous traffic



RFC 3918 IPv4 1-to-N Multicast Performance

With one transmitter port and 255 subscriber ports, IPv4 multicast tests used a traffic pattern similar to those in the Ethernet multicast tests. This time, however, the Cisco Nexus 9508 moved traffic across subnet boundaries in a test involving more than 1 million unique destinations .

In this Layer-3 test, engineers configured each interface on the Cisco switch in a different IPv4 subnet. For IPv4 multicast support, the switch ran the Protocol Independent Multicast-Sparse Mode (PIM) routing protocol as well as IGMPv3 to maintain IGMP snooping tables.

On the control plane, 255 subscriber ports each joined the same 4,095 IP multicast groups. On the data plane, the Spirent traffic generator offered traffic to all IP multicast groups on all 255 receiver interfaces. In total, then, the Cisco Nexus 9508 forwarded multicast traffic to 4,095 mroutes times 255 outgoing interfaces (OIFs), for a total of 1,044,225 unique destinations. Note that the use of 4,095 multicast groups was due to a stream count limit in the Spirent dX2 test modules, and is not a limit of the Cisco Nexus 9508.

The Cisco Nexus 9508 again delivered all traffic for all frame sizes with zero frame loss with throughput equivalent to 99.603 percent of line rate.

Table 14 presents throughput, latency, and jitter results for all frame sizes.

Figures 19 and 20 compare average and maximum delay measurements, respectively, with offered loads of 10, 50, 90, and 99.603 percent of line rate.

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	37,944,151,999	19.427	99.994%	1.250	3.138	3.480	0.010	0.080
128	21,535,870,055	22.053	99.994%	1.250	3.147	4.010	0.007	0.680
256	11,548,220,171	23.651	99.994%	1.260	3.178	3.520	0.006	0.070
512	5,991,181,893	24.540	99.994%	1.260	3.207	4.130	0.004	0.640
1,024	3,052,977,747	25.010	99.994%	1.280	3.261	3.570	0.004	0.040
1,280	2,451,775,976	25.106	99.994%	1.280	3.267	3.580	0.005	0.040
1,518	2,072,372,412	25.167	99.994%	1.280	3.273	4.400	0.004	0.620
9,216	345,096,231	25.443	99.994%	1.280	3.501	3.820	0.005	0.040

Table 14: IPv4 1-to-N multicast performance results

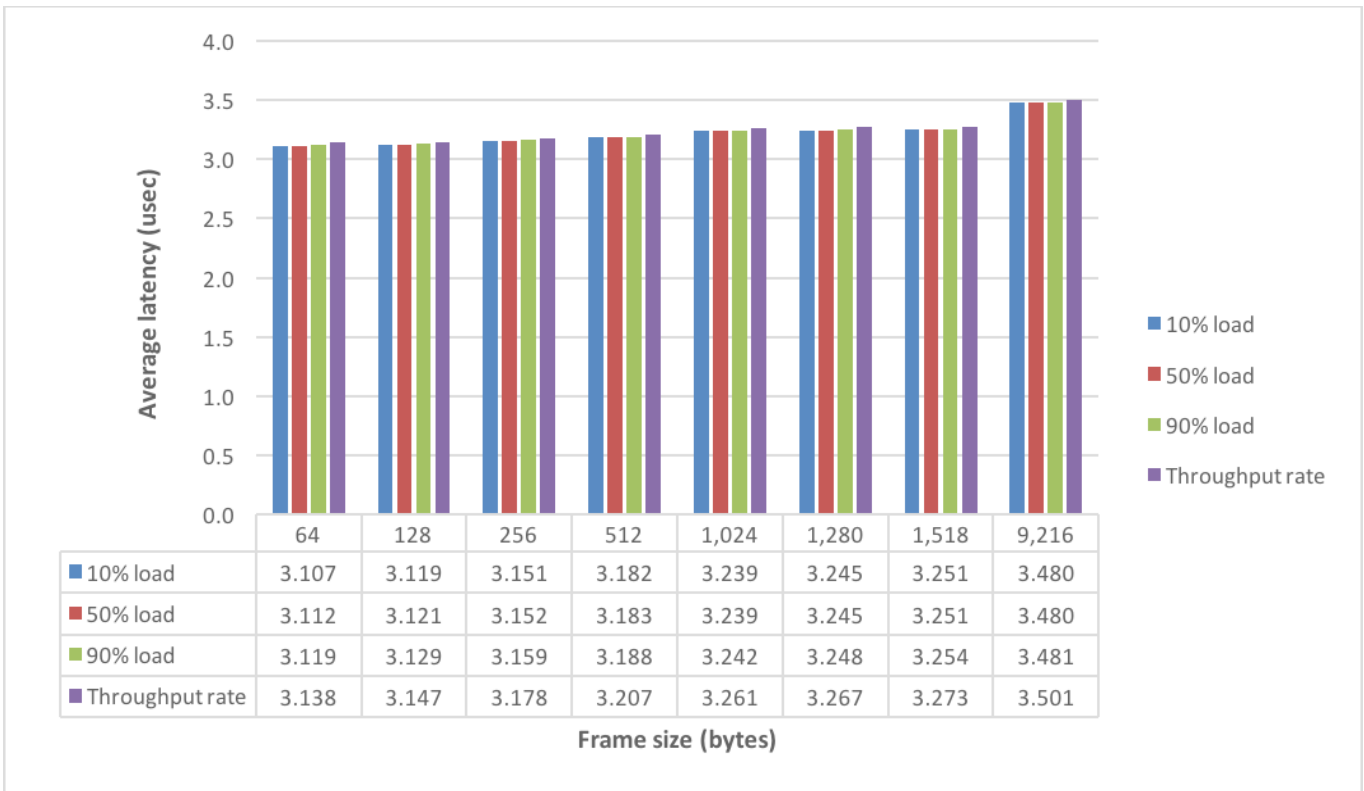


Figure 19: IPv4 1-to-N multicast average latency vs. load

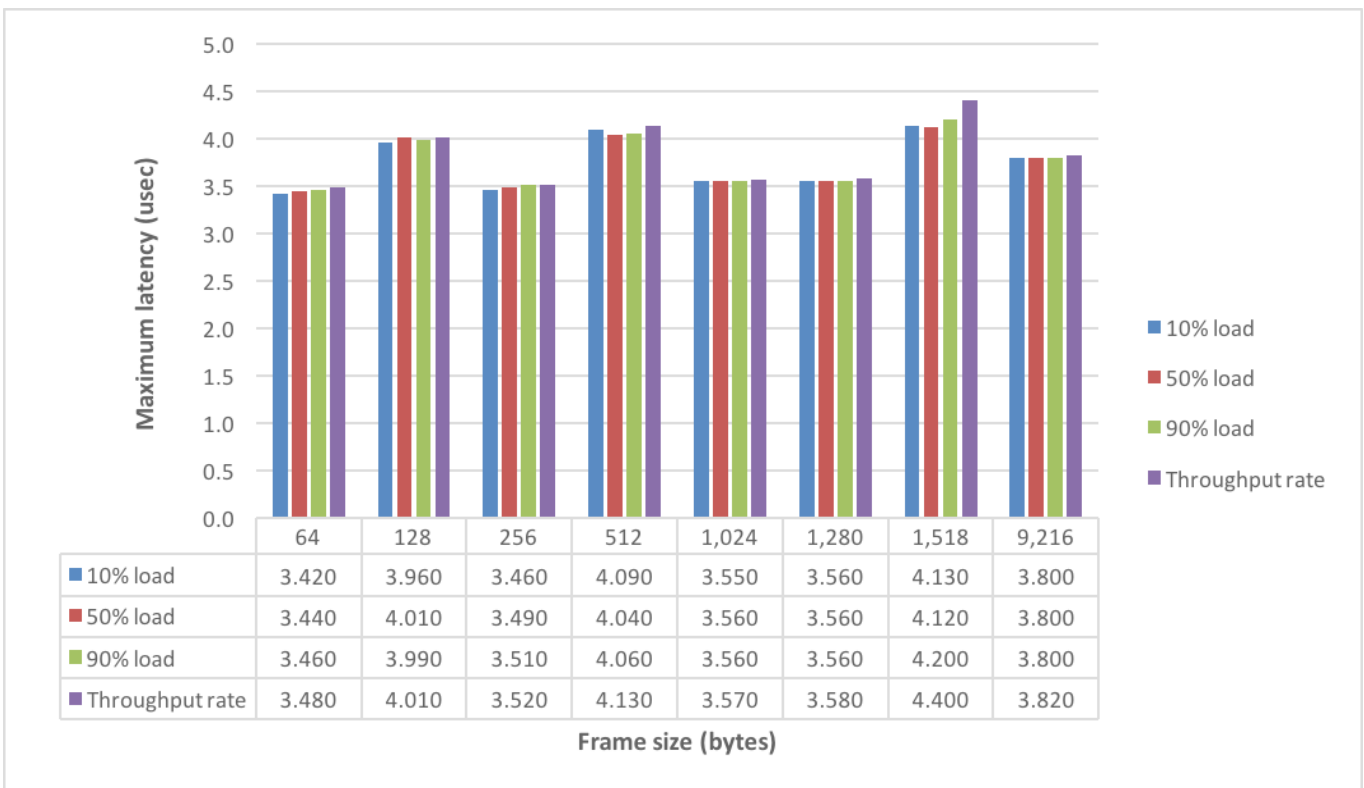


Figure 20: IPv4 1-to-N multicast maximum latency vs. load



IPv4 N-to-N Multicast Performance

Using a traffic pattern similar to the Ethernet case, the Layer-3 N-to-N tests involved simultaneous multicast transmit and receive ports, only this time the switch routed traffic to different subnets on each port. (See the “Ethernet N-to-N Multicast Performance” and “Test Methodology” sections for a complete description of the traffic pattern.)

In this scenario, engineers configured each port of the Spirent TestCenter instrument to offer traffic to one multicast group address, with the remaining 255 ports subscribed to that multicast group. With all ports configured this way, each Spirent port offered traffic to one multicast group address, and expected to receive traffic from the 255 other multicast group addresses on 255 other ports.

By definition, this was an overload test, since each output port received 255 frames at the same instant. In switch modules with relatively small buffers, the output ports may not have sufficient capacity to handle a 255:1 overload.

Fortunately, the new N9K-X9732C-EX modules for the Cisco Nexus 9508 have higher buffer capacity than previous products, and can handle this extreme load without frame loss. The tradeoff: Latency will naturally increase because of the larger buffers.

As in the 1-to-N multicast tests, the Cisco Nexus 9508 delivered N-to-N multicast traffic with zero frame loss in all test cases. Tables 15 and 16 present throughput, latency, and jitter results for all frame sizes with “synchronous start” and “asynchronous start” traffic, respectively.

Note that latency and jitter with staggered-start traffic – the only pattern that merchant silicon ASICs could handle – is far lower than with synchronous traffic. This is due to smaller buffers in merchant silicon ASICs.

Figures 21 and 22 compare average and maximum delay measurements with a “synchronous start,” using offered loads of 10, 50, 90, and the throughput rate.

Figures 23 and 24 compare average and maximum delay measurements with a “staggered start” (traffic on different ports begins asynchronously at 448-usec intervals), using offered loads of 10, 50, 90, and the throughput rate. .



Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,092,629,360	19.503	99.993%	1.030	3.518	4.710	0.090	1.220
128	21,620,141,008	22.139	99.993%	1.040	4.129	5.910	0.086	1.650
256	11,593,408,957	23.743	99.993%	1.030	5.388	8.480	0.134	3.160
512	6,014,625,713	24.636	99.993%	1.080	7.877	13.560	0.227	6.100
1,024	3,064,924,229	25.108	99.993%	1.120	12.848	23.740	0.437	11.950
1,280	2,461,369,926	25.204	99.993%	1.060	15.313	28.800	0.512	20.610
1,518	2,080,481,737	25.265	99.993%	1.070	17.603	33.520	0.960	25.730
9,216	346,446,621	25.543	99.993%	1.040	95.158	189.120	10.477	181.040

Table 15: IPv4 N-to-N multicast performance results with synchronous traffic

Frame size (bytes)	Throughput			Latency			Jitter	
	Frames/s	Tbit/s	% line rate	Min (usec)	Avg (usec)	Max (usec)	Avg (usec)	Max (usec)
64	38,020,239,903	19.466	99.803%	1.000	2.919	3.790	0.008	0.110
128	21,579,055,542	22.097	99.803%	1.010	2.943	4.380	0.007	0.710
256	11,571,378,111	23.698	99.803%	1.020	3.070	4.100	0.014	0.710
512	6,003,196,681	24.589	99.803%	1.030	3.091	3.990	0.009	0.240
1,024	3,059,100,756	25.060	99.803%	1.050	3.425	4.500	0.007	0.280
1,280	2,456,693,435	25.157	99.803%	1.050	3.265	4.560	0.008	0.760
1,518	2,076,529,079	25.217	99.803%	1.060	3.656	4.900	0.035	1.020
9,216	345,789,310	25.494	99.803%	1.050	4.570	7.150	0.009	1.460

Table 16: IPv4 N-to-N multicast performance results with asynchronous traffic

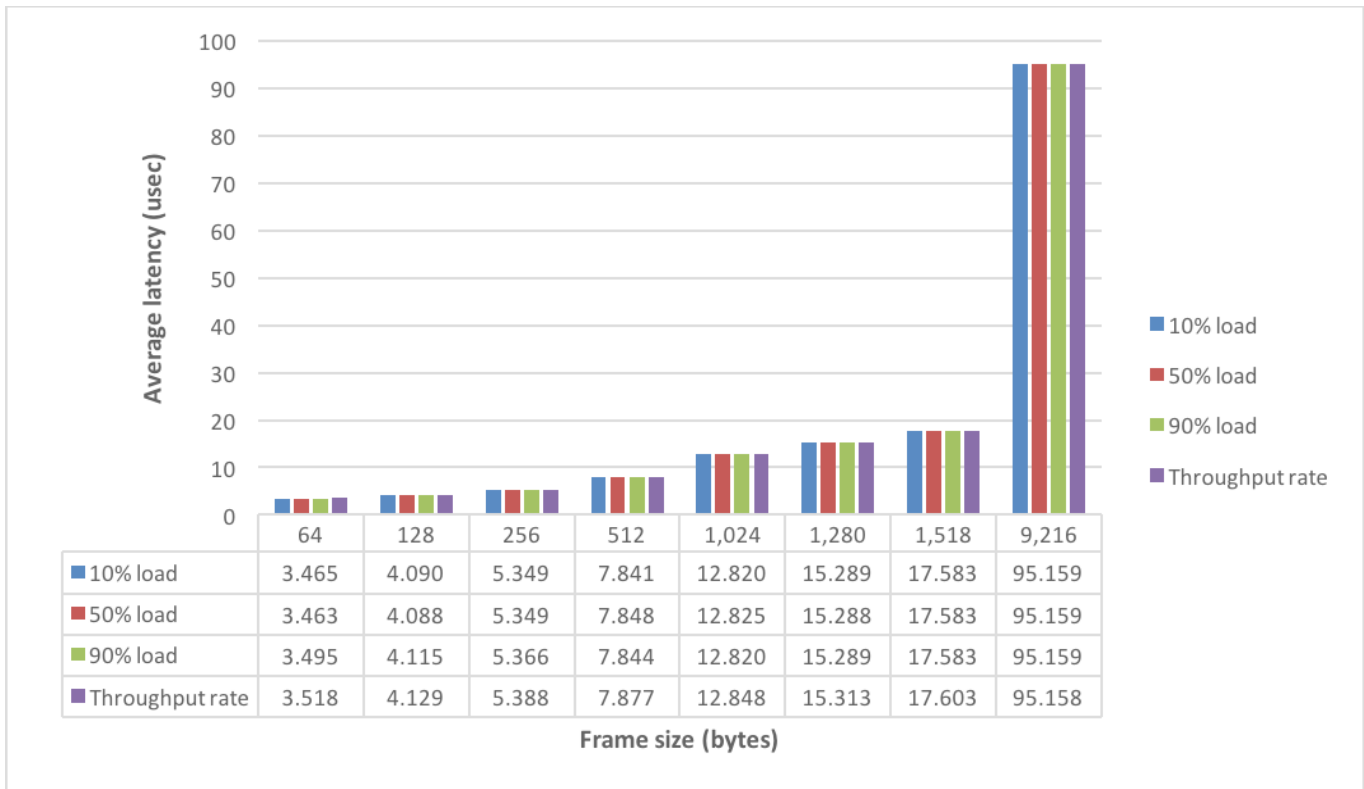


Figure 21: IPv4 N-to-N multicast average latency vs. load with synchronous traffic

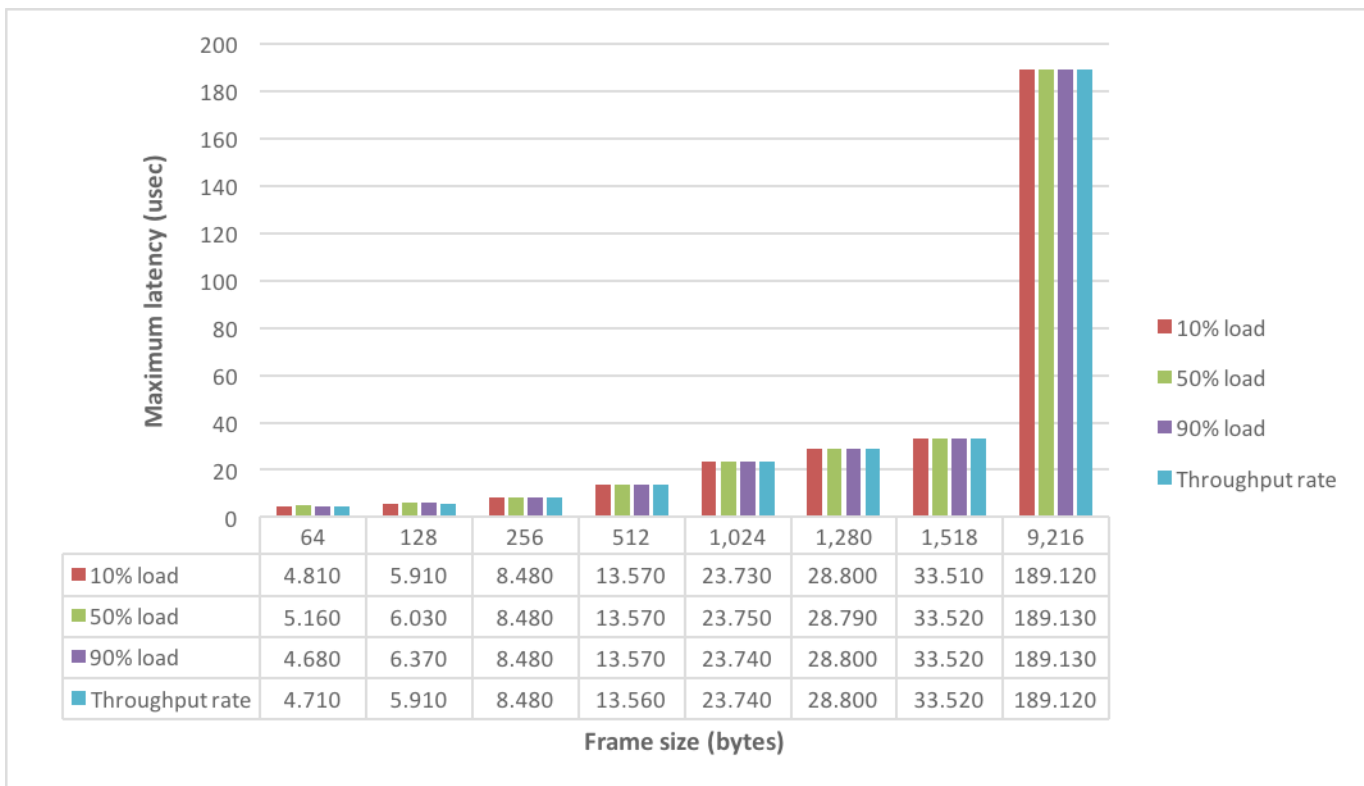


Figure 22: IPv4 N-to-N multicast maximum latency vs. load with synchronous traffic

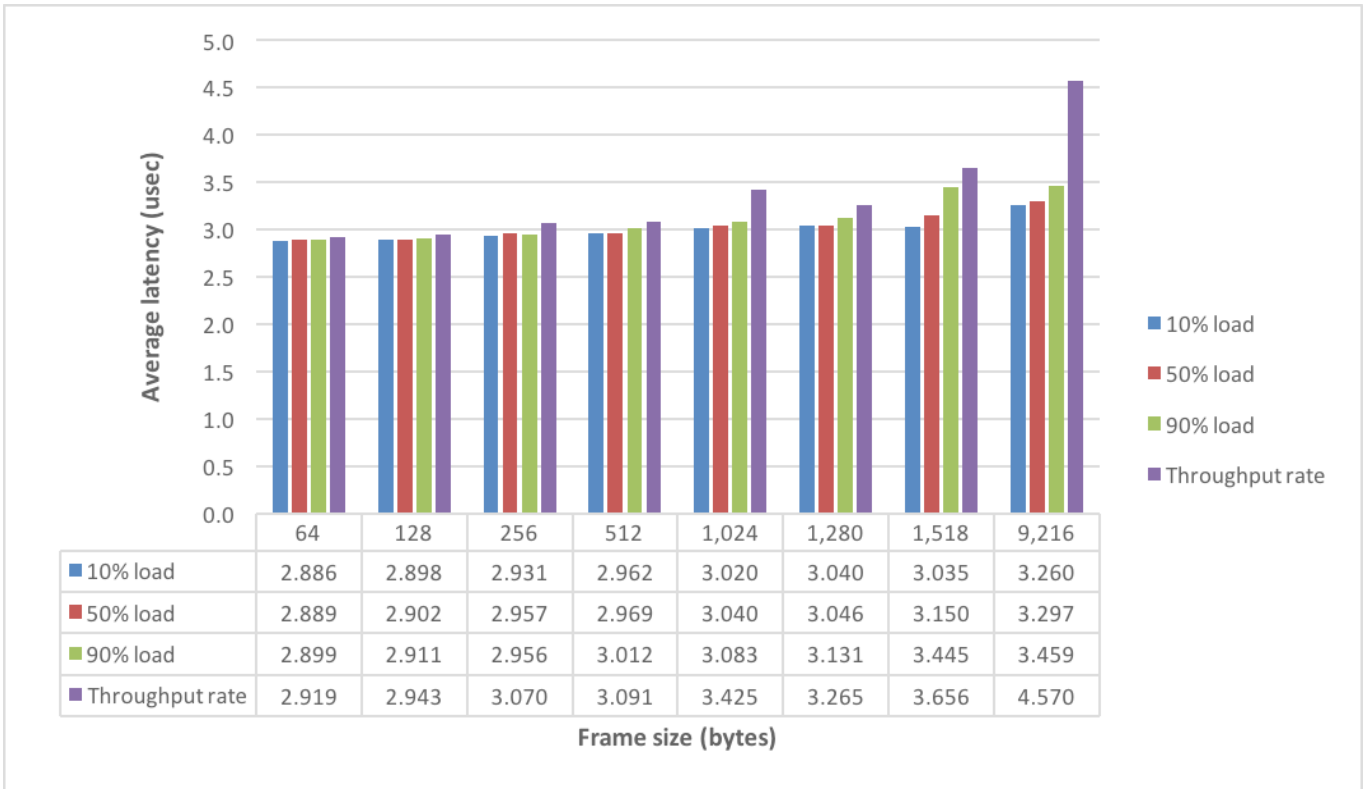


Figure 23: IPv4 N-to-N multicast average latency vs. load with asynchronous traffic

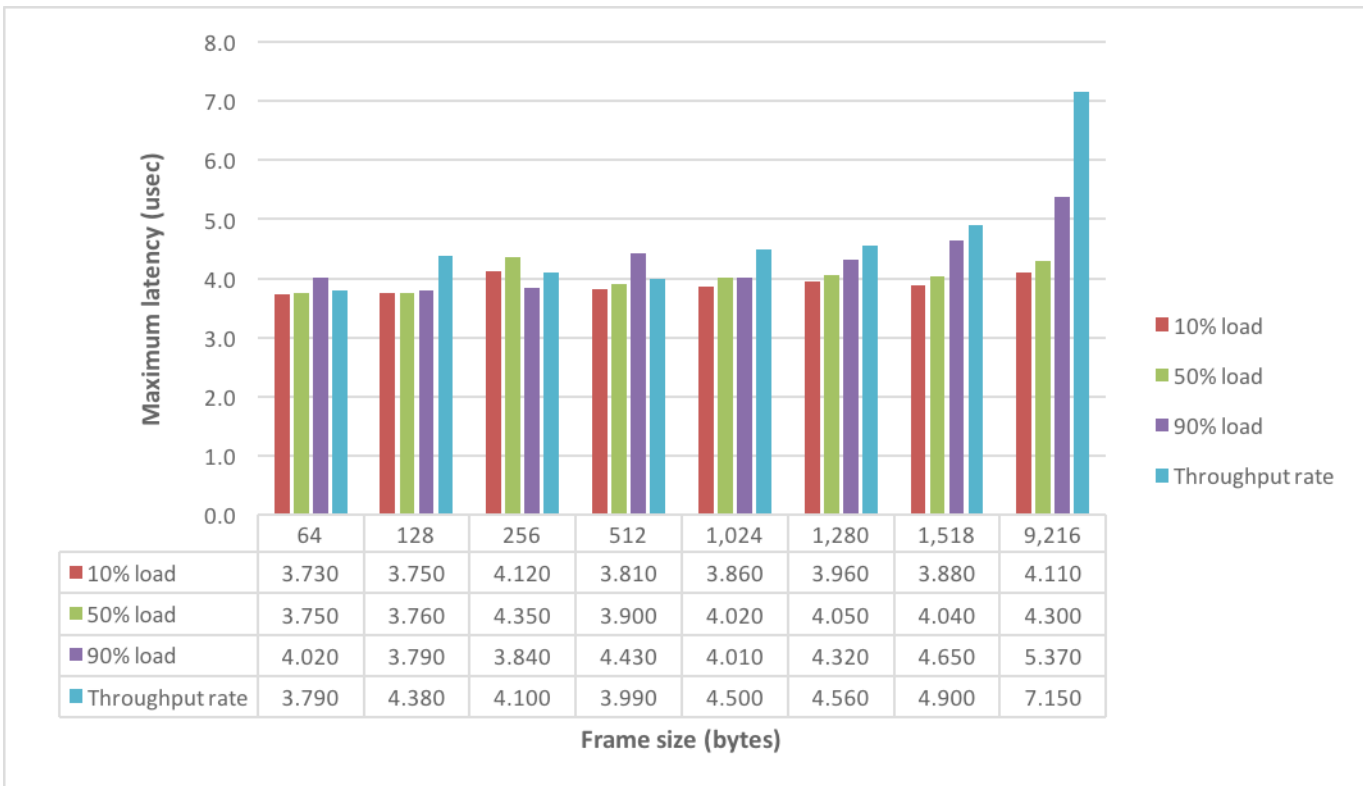


Figure 24: IPv4 N-to-N multicast maximum latency vs. load with asynchronous traffic



Power Consumption

For both ecological and financial reasons, operators of data centers face more pressure than ever to reduce power consumption. While servers continue to represent the greatest power cost in most data centers, core switches also make a significant contribution. Increasingly, customers include power usage among their selection criteria when evaluating new data center equipment.

Engineers measured power consumption in two modes:

- Switch idle with all 256 100G Ethernet transceivers in place (link up)
- Switch forwarding unicast frames at the throughput rate on all 256 100G Ethernet ports in a fully meshed pattern, using each of the frame sizes from the unicast throughput tests

Engineers used a Fluke clamp meter to measure line voltage at the power source and a Voltech Power Analyzer to measure amperage at the power supply, and then calculated watts by multiplying volts and amps.

The Cisco Nexus 9508 has eight power supplies, with four apiece arrayed in two grids for redundancy. At any one time, the system uses one grid and load-shares current across the power supplies in that grid. After verifying load-sharing was approximately equal across power supplies, engineers measured power usage on two supplies and then multiplied by 2 to obtain total wattage.

In the worst-case scenario, the Cisco Nexus 9508 consumed about 37 watts per port when forwarding 64-byte frames on all ports. Power consumption falls dramatically in test cases with longer frames, and with no traffic. When idle, power usage drops to about 21 watts per port, and rises only to 23 watts per port when forwarding jumbo frames on all ports in a fully meshed pattern.

Significantly, the power-per-port measurements presented here are derived from total system wattage divided by port count. When assessing power consumption, it is important to consider whether wattage numbers take into account power draw for transceivers, a switch chassis, fan, fabric, and other components. *The power consumption numbers presented here take all components, including transceivers, into account.*

Table 17 presents total system power consumption measurements for the 100G Ethernet test bed when idle, and with all frame sizes used in the throughput and latency tests.

Interface count (with transceivers)	Frame size (bytes)	Watts per port	Total watts
256 x 100G	NA (idle)	20.91	5,354.00
256 x 100G	64	36.99	9,470.60
256 x 100G	128	31.13	7,969.93
256 x 100G	256	27.91	7,144.71
256 x 100G	512	25.59	6,552.12
256 x 100G	1,024	24.04	6,155.05
256 x 100G	1,280	24.00	6,143.83
256 x 100G	1,518	23.90	6,118.79
256 x 100G	9,216	23.36	5,979.82

Table 17: Cisco 9508 power consumption



Forward Error Correction (FEC) Latency and Jitter

Ethernet transceivers use a mechanism called Forward Error Correction (FEC) to detect and correct errors without the need for retransmission. Most Ethernet devices have FEC enabled by default, with users willing to trade off a small amount of added latency for protection against low-level errors.

However, when applications require the absolute lowest latency, network professionals may disable FEC toward that end. Although Cisco recommends leaving FEC enabled, there are data-center applications, particularly those involving fiber-optic transceivers and cable lengths of 3 meters or less, where users may opt to disable FEC.

To determine the effect of FEC on latency, engineers offered fully meshed traffic to a single 32-port Cisco N9K-X9732C-EX module with copper cabling, both with FEC enabled and disabled.

Table 18 presents the differences between latencies with FEC enabled and disabled. As the results show, disabling FEC can reduce latency by more than 260 nanoseconds, on average, for users willing to forgo forward error correction. Minimum latency is omitted here because the minimum delta was 0 in all cases.

Also, sharp-eyed readers may notice that maximum values are less than average values for some frame sizes. Since the table presents differences in measurements rather than the measurements themselves, it is possible to have cases where average deltas (differences) exceed maximum deltas.

These tests involved fully meshed traffic, the most stressful traffic pattern, with 32 ports. For the sake of completeness, engineers repeated the 2-port tests with FEC disabled (see the “Testing Two Ports at a Time” section). Tables 19, 20 and 21 present the differences between latencies with FEC enabled and disabled in 2-port tests in same-ASIC, ASIC-to-ASIC, and module-to-module test cases. Even with just 2 ports, disabling FEC can reduce latency by up to 260 nanoseconds, on average. Minimum latency again is omitted because the minimum delta was 0 in all cases.

Frame size (bytes)	Latency			Jitter	
	Minimum (usec)	Average (usec)	Maximum (usec)	Average (usec)	Maximum (usec)
64	0.260	0.268	-0.370	0.000	-0.650
128	0.250	0.271	0.290	0.001	0.020
256	0.250	0.270	0.270	0.001	0.010
512	0.260	0.267	-0.500	0.001	-0.800
1,024	0.260	0.268	0.270	0.000	0.000
1,280	0.250	0.268	1.340	-0.002	1.100
1,518	0.270	0.267	-0.650	0.001	-1.000
9,216	0.260	0.268	0.270	-0.001	0.010

Table 18: Difference in latency with FEC enabled and disabled, 32-port test



Frame size (bytes)	Latency			Jitter	
	Minimum (usec)	Average (usec)	Maximum (usec)	Average (usec)	Maximum (usec)
64	0.230	0.234	0.230	0.001	0.000
128	0.230	0.232	0.240	0.000	0.000
256	0.240	0.235	0.230	0.000	0.000
512	0.240	0.233	0.230	0.000	0.000
1,024	0.240	0.233	0.230	0.000	0.000
1,280	0.240	0.233	0.230	0.000	0.000
1,518	0.240	0.233	0.240	0.000	0.000
9,216	0.230	0.234	0.230	0.000	0.010

Table 19: Difference in latency with FEC enabled and disabled, same ASIC

Frame size (bytes)	Latency			Jitter	
	Minimum (usec)	Average (usec)	Maximum (usec)	Average (usec)	Maximum (usec)
64	0.240	0.239	0.240	0.000	0.000
128	0.240	0.237	0.240	0.000	0.000
256	0.240	0.238	0.230	0.000	0.000
512	0.240	0.238	0.230	0.000	0.010
1,024	0.230	0.238	0.250	-0.001	0.010
1,280	0.230	0.237	0.250	-0.001	0.010
1,518	0.230	0.239	0.240	0.000	0.000
9,216	0.230	0.239	0.240	0.000	0.000

Table 20: Difference in latency with FEC enabled and disabled, same module, different ASICs



Frame size (bytes)	Latency			Jitter	
	Minimum (usec)	Average (usec)	Maximum (usec)	Average (usec)	Maximum (usec)
64	0.250	0.255	0.240	0.000	0.000
128	0.260	0.253	0.250	0.000	0.000
256	0.260	0.254	0.250	0.000	0.000
512	0.260	0.254	0.240	0.000	0.000
1,024	0.260	0.253	0.250	-0.001	0.000
1,280	0.270	0.253	0.250	-0.001	0.000
1,518	0.260	0.253	0.240	-0.001	0.000
9,216	0.260	0.253	0.240	0.000	0.000

Table 21: Difference in latency with FEC enabled and disabled, same module, different modules

Test Methodology

The principle objective of this test was to characterize the performance of the Cisco Nexus 9508 equipped with 256 100G Ethernet interfaces in various Layer-2 and Layer-3 configurations. Network Test evaluated the Cisco Nexus 9508 in 13 scenarios, plus one additional case with no Cisco switch present:

- Test bed infrastructure latency and jitter
- Port-to-port performance
- Cisco Cloudscale Technology vs. merchant silicon throughput
- RFC 2889 Ethernet unicast performance
- RFC 2544 IPv4 unicast performance
- RFC 2544 IPv4 unicast performance with BGP routing
- RFC 5180 IPv6 unicast performance
- RFC 5180 IPv6 unicast performance with BGP-MP routing
- RFC 3918 Ethernet 1-to-N multicast performance
- Ethernet N-to-N multicast performance
- RFC 3918 IPv4 1-to-N multicast performance
- IPv4 N-to-N Multicast Performance
- Power consumption
- Forward error correction (FEC) latency and jitter

For all configurations, the performance metrics consisted of throughput; minimum, average, and maximum latency; and average and maximum jitter. The test results presented here omit minimum jitter because it was 0 in all test cases.



The principle test instrument for this project was the Spirent TestCenter traffic generator/analyzer equipped with dX2-100G-P4 modules. For unicast tests, the Spirent instrument offered traffic to all 256 100G Ethernet ports in a fully meshed pattern, meaning all traffic was destined for all other ports. For RFC 3918 multicast tests, the Spirent instrument used the IGMPv3 protocol to subscribe to 4,095 IP multicast group addresses on 255 receiver ports. A single Spirent transmitter port then offered traffic to all IP multicast group addresses in the Ethernet and IPv4 multicast test cases. In the N-to-N multicast tests, all 256 ports were simultaneously multicast transmitters and receivers.

In all unicast test cases, the Spirent test instrument offered traffic at a maximum of 99.994 percent of line rate for a duration of 60 seconds, and measured the latency of every frame received. Previous Network Test assessments of Cisco switches used 300-second durations, in part because of customer interest and in part because earlier switches had smaller buffers that could cause latency and jitter to rise over time. Testing with the N9K-X9732C-EX showed no significant difference in latency or jitter between 60- and 300-second test durations.

In unicast test cases, engineers used 99.994 percent of line rate, which is 60 parts per million (60 ppm) slower than nominal line rate, to avoid clocking differences between the traffic generator and the switch under test. The IEEE 802.3 Ethernet specification requires interfaces to tolerate clocking differences of up to +/- 100 ppm, so a 60-ppm difference is well within that specification.

Test engineers repeated this test with eight frame sizes: 64-, 128-, 256-, 512-, 1,024-, 1,280, 1,518-, and 9,216-byte Ethernet frames. The first seven sizes are recommended in RFC 2544, while data-center applications that involve high-volume data transfer often use 9,216-byte jumbo frames.

Engineers configured the Spirent instrument to measure latency using the first-in, first-out (FIFO) measurement method described in RFC 1242. FIFO latency measurement is appropriate when switches are configured in so-called cut-through mode, where the switch begins forwarding each incoming frame immediately, before caching the entire frame. The Cisco Nexus 9508 also can be configured to operate in store-and-forward mode, where the switch caches the entire frame before forwarding each it, albeit with proportionately higher latency as frame size increases.

Engineers also measured switch delay for three loads lower than the throughput rate – at 10, 50, and 90 percent of line rate. RFC 2544 requires latency to be measured at, and only at, the throughput rate. Since production networks typically see far lower average utilization, Cisco requested additional tests to be run to characterize delay at lower offered loads.

IPv6 tests used 86-byte instead of 64-byte frames as the minimum frame length. This is due to two requirements. First, the Spirent test instrument embeds a 20-byte “signature field” in every test frame. Second, in most tests test engineers configured traffic to use an 8-byte UDP header to compare results with earlier tests using merchant silicon ASICs; in that earlier project, a large amount of UDP randomness was needed to ensure optimal distribution of flows across the internal switch fabric. Adding up all the field lengths (18 bytes for Ethernet header and CRC; 40 bytes for IPv6 header; 8 bytes for UDP header; and 20 bytes for Spirent signature field) yields a minimum frame size of 86 bytes.

Some IPv4 and IPv6 unicast tests involved direct routes (1 per port), while others used Border Gateway Protocol (BGP). In both IPv4 and IPv6 tests, each Spirent test interface represented one BGP router advertising reachability to 8 networks, for a total of 2,048 unique networks. Both IPv4 and IPv6 network counts represent limits of the Spirent dX2 test modules and not that of the Cisco Nexus 9508.

The RFC 3918 Ethernet multicast traffic tests involved a traffic pattern with one transmitter port and 255 receiver (subscriber) ports. Here, all 255 receiver ports on the Spirent TestCenter instrument joined the same 4,095 multicast groups using IGMPv3 reports. After the switch’s IGMP snooping table was fully populated, the test



instrument then offered traffic to the single transmit port, with destination addresses of all 4,095 multicast groups. As in the unicast tests, the instrument measured throughput and latency for eight frame sizes. This and all other multicast tests used group addresses beginning at 225.0.1.0/32 and incrementing by 1.

The N-to-N multicast tests also involved a single VLAN, but this time engineers configured all 256 ports to be multicast transmitters. In this scenario, engineers configured Spirent TestCenter to offer multicast traffic on each port destined to a different multicast group address. Thus, in all, there were 256 multicast transmitters, each sending to one unique multicast group address, and 255 ports receiving traffic from all group addresses other than that of their own transmitter. Engineers repeated the N-to-N multicast tests with the switch in Layer-2 and Layer-3 configurations.

The layer-3 RFC 3918 multicast tests used the same traffic pattern as the layer-2 tests, with one transmitter port and 255 receiver (subscriber) ports. In this case, however, all switch ports also ran the protocol independent multicast-sparse mode (PIM-SM) routing protocol. All switch ports used PIM-SM to learn multicast routes. Then, all 255 receiver ports on the Spirent TestCenter instrument joined the same 4,095 multicast groups using IGMPv3 reports. The instrument measured throughput and latency for the same eight frame sizes as in the other performance tests.

Notably, test engineers did not configure the Spirent test instrument with latency compensation or parts-per-million (PPM) clocking adjustments. These adjustments exist in test instruments to compensate for very specific use cases, but also can be abused. The adjustment of time measurements in a test instrument for purposes of “improving” test results is generally considered to be an unscrupulous practice.

For reproducibility of these results, it’s important to note the contents of test traffic, especially with regard to MAC and IP addresses and UDP port numbers. In the Layer-2 unicast tests, all Spirent emulated hosts used pseudorandom MAC addresses as described in RFC 4814. The Spirent IP addresses began at 10.0.0.2/16, incrementing by port number. All Cisco Nexus 9508 interfaces were members of the same VLAN, which was bound to an IPv4 address of 10.0.0.1/8 (though this was not used in this Layer-2 test). The UDP headers used 8,000 unique source and destination ports, each beginning at 20001 and incrementing by 1 up to 28,000. In tests involving BGP and BGP-MP, UDP headers used random source and destination port numbers. The Layer-2 multicast tests used Spirent default MAC addresses and IPv4 addresses starting at 10.0.0.2/16 and incrementing by port number.

Conclusion

These results set new speed records in data-center switching, with the highest throughput and lowest latency and jitter ever recorded on a 256-port 100G Ethernet test bed. The Cisco Nexus 9508 never dropped a frame in rigorous benchmarks covering unicast, multicast, Ethernet, IPv4, IPv6, BGP traffic, all with traffic moving at virtual line rate across its 256 100G Ethernet ports. Moreover, tests showed the Nexus 9508 to be a highly capable performer both in 1-to-N and N-to-N multicast scenarios.

Latency and jitter also remained low and constant across test cases, a critical requirement for time-sensitive applications. Average and maximum delay is lower still in test cases involving traffic at 10, 50, and 90 percent of wire speed, providing a complete picture of how the switch is likely to perform in production settings. And customers who require still lower latency and jitter for some applications can achieve reductions of more than 250 ns by disabling forward error correction.

For network professionals looking to build the very largest data centers, and for those just looking to ensure a pathway for future growth, the Cisco Nexus 9508 with Cisco Cloudscale Technology proved highly capable across all these rigorous tests.



Appendix A: Jitter Measurements

This section presents average and maximum jitter measurements across varying offered loads, ranging from 10 percent to the throughput rate. Jitter, or delay variation, is a critical metric for any application sensitive to delay. High jitter can severely degrade the performance of video and video applications, as well as any other application that requires real-time delivery of messages.

As mentioned earlier, minimum jitter measurements were 0 and thus are omitted here. This appendix also omits average and maximum jitter measurements at various loads from the 2-port tests because they often are below the measurement resolution of the test instrument, and thus not meaningful. See the “Testing Two Ports at a Time” section for average and maximum jitter measurements at the throughput rate.

Table 22 presents average jitter measurements from the RFC 2889 Ethernet unicast performance tests.

Frame size (bytes)	Average jitter (usec)			
	10% load	50% load	90% load	Throughput rate
64	0.007	0.007	0.010	0.009
128	0.007	0.007	0.009	0.007
256	0.007	0.006	0.008	0.007
512	0.007	0.007	0.008	0.007
1,024	0.008	0.007	0.008	0.005
1,280	0.007	0.007	0.008	0.005
1,518	0.007	0.007	0.007	0.006
9,216	0.007	0.007	0.008	0.008

Table 22: Ethernet unicast average jitter

Table 23 presents maximum jitter measurements from the RFC 2889 Ethernet unicast performance tests.

Frame size (bytes)	Maximum jitter (usec)			
	10% load	50% load	90% load	Throughput rate
64	0.700	0.760	1.080	1.140
128	0.060	0.110	0.330	0.430
256	0.070	0.130	0.300	0.500
512	0.690	0.740	1.120	1.180
1,024	0.130	0.230	0.830	1.150
1,280	0.690	0.770	1.400	1.780
1,518	0.170	0.180	1.090	1.410
9,216	0.670	0.780	1.420	1.970

Table 23: Ethernet unicast maximum jitter



Table 24 presents average jitter measurements from the RFC 2544 IPv4 unicast performance tests.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.007	0.007	0.010	0.009
128	0.007	0.007	0.009	0.008
256	0.007	0.006	0.009	0.007
512	0.007	0.007	0.008	0.007
1,024	0.008	0.007	0.008	0.004
1,280	0.007	0.007	0.008	0.005
1,518	0.007	0.007	0.007	0.006
9,216	0.007	0.007	0.008	0.007

Table 24: IPv4 unicast average jitter

Table 25 presents maximum jitter measurements from the RFC 2544 IPv4 unicast performance tests.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.060	0.120	0.950	1.140
128	0.060	0.700	0.310	0.490
256	0.700	0.140	0.930	0.550
512	0.060	0.720	0.580	1.070
1,024	0.700	0.240	1.490	1.180
1,280	0.070	0.760	0.930	1.880
1,518	0.700	0.170	1.140	1.450
9,216	0.060	0.780	1.440	2.010

Table 25: IPv4 unicast maximum jitter



Table 26 presents average jitter measurements from the RFC 2544 IPv4 unicast performance tests with BGP routing enabled.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.007	0.007	0.007	0.006
128	0.007	0.007	0.006	0.007
256	0.007	0.007	0.006	0.006
512	0.007	0.008	0.007	0.006
1,024	0.007	0.007	0.008	0.007
1,280	0.006	0.007	0.008	0.006
1,518	0.007	0.008	0.007	0.007
9,216	0.007	0.007	0.008	0.007

Table 26: IPv4/BGP unicast average jitter

Table 27 presents maximum jitter measurements from the RFC 2544 IPv4 unicast performance tests with BGP routing enabled.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.060	0.110	1.490	1.650
128	0.700	0.140	0.890	1.280
256	0.070	0.720	1.410	1.800
512	0.060	0.120	1.500	2.080
1,024	0.700	0.130	1.480	2.130
1,280	0.060	0.760	1.490	2.070
1,518	0.090	0.170	1.490	2.100
9,216	0.490	0.770	1.420	1.980

Table 27: IPv4/BGP unicast maximum jitter



Table 28 presents average jitter measurements from the RFC 5180 IPv6 unicast performance tests.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
86	0.007	0.007	0.009	0.009
128	0.007	0.007	0.009	0.008
256	0.007	0.006	0.009	0.007
512	0.007	0.007	0.009	0.007
1,024	0.008	0.007	0.008	0.004
1,280	0.007	0.007	0.008	0.005
1,518	0.007	0.007	0.008	0.006
9,216	0.007	0.007	0.008	0.008

Table 28: IPv6 unicast average jitter

Table 29 presents maximum jitter measurements from the RFC 5180 IPv6 unicast performance tests.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
86	0.710	0.710	0.380	0.520
128	0.060	0.110	1.100	1.090
256	0.060	0.730	0.340	0.550
512	0.680	0.130	0.410	0.640
1,024	0.070	0.140	1.150	1.170
1,280	0.700	0.760	0.390	0.630
1,518	0.080	0.180	1.210	1.370
9,216	0.060	0.680	0.240	0.480

Table 29: IPv6 unicast maximum jitter



Table 30 presents average jitter measurements from the RFC 5180 IPv6 unicast performance tests with BGP-MP routing enabled.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
86	0.007	0.006	0.007	0.007
128	0.007	0.007	0.006	0.006
256	0.007	0.007	0.007	0.006
512	0.007	0.007	0.007	0.006
1,024	0.007	0.008	0.007	0.006
1,280	0.007	0.008	0.007	0.006
1,518	0.007	0.008	0.008	0.007
9,216	0.007	0.008	0.008	0.007

Table 30: IPv6/BGP unicast average jitter

Table 31 presents maximum jitter measurements from the RFC 5180 IPv6 unicast performance tests with BGP-MP routing enabled.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
86	0.060	0.080	1.350	0.890
128	0.060	0.730	0.440	0.700
256	0.700	0.090	1.170	0.660
512	0.060	0.740	0.400	1.350
1,024	0.700	0.140	1.350	0.590
1,280	0.070	0.780	0.310	1.700
1,518	0.690	0.180	0.340	0.580
9,216	0.060	0.070	0.700	0.650

Table 31: IPv6/BGP unicast maximum jitter



Table 32 presents average jitter measurements from the RFC 3918 Ethernet 1-to-N multicast tests.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.009	0.008	0.007	0.007
128	0.009	0.008	0.006	0.005
256	0.009	0.009	0.006	0.004
512	0.009	0.010	0.006	0.004
1,024	0.010	0.009	0.008	0.004
1,280	0.009	0.009	0.008	0.004
1,518	0.009	0.009	0.008	0.004
9,216	0.009	0.009	0.009	0.004

Table 32: Ethernet 1-to-N multicast average jitter

Table 33 presents maximum jitter measurements from the RFC 3918 Ethernet 1-to-N multicast tests.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.600	0.690	0.680	0.700
128	0.080	0.070	0.060	0.070
256	0.070	0.080	0.060	0.060
512	0.620	0.650	0.680	0.650
1,024	0.070	0.080	0.070	0.040
1,280	0.080	0.070	0.070	0.040
1,518	0.320	0.620	0.570	0.620
9,216	0.070	0.080	0.080	0.050

Table 33: Ethernet 1-to-N multicast maximum jitter



Table 34 presents average jitter measurements from the Ethernet N-to-N multicast tests with synchronous traffic. Note that results in Tables 34 and 35 are the result of an overload case, with 255 frames presented to each of 256 egress interfaces at the same instant.

	Average jitter (usec)			
Frame size (bytes)	10% load	50% load	90% load	Throughput rate
64	0.074	0.063	0.070	0.115
128	0.128	0.104	0.141	0.104
256	0.250	0.176	0.224	0.160
512	0.441	0.291	0.504	0.282
1,024	0.762	0.559	0.826	0.566
1,280	0.906	0.654	1.261	0.662
1,518	1.589	1.448	0.971	1.215
9,216	13.137	14.128	11.799	13.840

Table 34: Ethernet N-to-N multicast average jitter with synchronous traffic

Table 35 presents maximum jitter measurements from the Ethernet N-to-N multicast tests with synchronous traffic.

	Maximum jitter (usec)			
Frame size (bytes)	10% load	50% load	90% load	Throughput rate
64	0.950	0.950	1.060	1.570
128	2.110	2.510	1.690	2.860
256	3.120	3.130	4.480	3.130
512	6.270	6.300	6.270	6.320
1,024	12.470	12.280	12.300	17.800
1,280	23.740	23.840	23.720	23.820
1,518	29.550	29.650	29.780	29.790
9,216	181.050	180.330	181.050	181.060

Table 35: Ethernet N-to-N multicast maximum jitter with synchronous traffic



Table 36 presents average jitter measurements from the Ethernet N-to-N multicast tests with asynchronous traffic.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.007	0.006	0.008	0.008
128	0.007	0.006	0.009	0.007
256	0.006	0.009	0.010	0.014
512	0.007	0.007	0.011	0.009
1,024	0.007	0.006	0.010	0.007
1,280	0.007	0.008	0.008	0.007
1,518	0.008	0.016	0.024	0.036
9,216	0.007	0.007	0.019	0.008

Table 36: Ethernet N-to-N multicast average jitter with asynchronous traffic

Table 37 presents maximum jitter measurements from the Ethernet N-to-N multicast tests with asynchronous traffic.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.060	0.710	0.110	0.100
128	0.070	0.100	0.110	0.110
256	0.680	0.170	0.150	0.690
512	0.620	0.750	0.800	0.270
1,024	0.130	0.210	0.800	0.280
1,280	0.160	0.250	0.350	0.720
1,518	0.590	0.860	0.800	1.490
9,216	0.140	1.300	1.090	1.300

Table 37: Ethernet N-to-N multicast maximum jitter with asynchronous traffic



Table 38 presents average jitter measurements from the RFC 3918 IPv4 1-to-N multicast tests.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.012	0.010	0.009	0.010
128	0.012	0.010	0.008	0.007
256	0.011	0.011	0.007	0.006
512	0.012	0.012	0.007	0.004
1,024	0.012	0.012	0.008	0.004
1,280	0.012	0.011	0.009	0.005
1,518	0.011	0.012	0.010	0.004
9,216	0.012	0.012	0.011	0.005

Table 38: IPv4 1-to-N multicast average jitter

Table 39 presents maximum jitter measurements from the RFC 3918 IPv4 1-to-N multicast tests.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.080	0.080	0.080	0.080
128	0.640	0.680	0.660	0.680
256	0.080	0.090	0.070	0.070
512	0.660	0.650	0.670	0.640
1,024	0.080	0.090	0.070	0.040
1,280	0.080	0.090	0.080	0.040
1,518	0.680	0.620	0.570	0.620
9,216	0.080	0.090	0.090	0.040

Table 39: IPv4 1-to-N multicast maximum jitter



Table 40 presents average jitter measurements from the IPv4 N-to-N multicast tests with synchronous traffic. Note that results in Tables 40 and 41 are the result of an overload case, with 255 frames presented to each of 256 egress interfaces at the same instant.

	Average jitter (usec)			
Frame size (bytes)	10% load	50% load	90% load	Throughput rate
64	0.065	0.052	0.065	0.090
128	0.110	0.086	0.115	0.086
256	0.195	0.141	0.179	0.134
512	0.343	0.239	0.347	0.227
1,024	0.618	0.447	0.616	0.437
1,280	0.729	0.521	0.882	0.512
1,518	1.214	1.131	0.778	0.960
9,216	10.933	9.391	8.221	10.477

Table 40: IPv4 N-to-N multicast average jitter with synchronous traffic

Table 41 presents maximum jitter measurements from the IPv4 N-to-N multicast tests with synchronous traffic.

	Maximum jitter (usec)			
Frame size (bytes)	10% load	50% load	90% load	Throughput rate
64	1.220	1.540	1.080	1.220
128	1.610	2.610	2.800	1.650
256	3.140	3.130	3.160	3.160
512	6.090	6.090	6.100	6.100
1,024	11.870	13.780	14.950	11.950
1,280	20.080	20.600	19.870	20.610
1,518	25.590	25.480	25.490	25.730
9,216	179.580	179.590	181.040	181.040

Table 41: IPv4 N-to-N multicast maximum jitter with synchronous traffic



Table 42 presents average jitter measurements from the IPv4 N-to-N multicast tests with asynchronous traffic.

Frame size (bytes)	Average jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.007	0.006	0.008	0.008
128	0.007	0.007	0.009	0.007
256	0.006	0.009	0.009	0.014
512	0.007	0.007	0.011	0.009
1,024	0.007	0.006	0.010	0.007
1,280	0.008	0.008	0.010	0.008
1,518	0.009	0.016	0.033	0.035
9,216	0.007	0.007	0.009	0.009

Table 42: IPv4 N-to-N multicast average jitter with asynchronous traffic

Table 43 presents maximum jitter measurements from the IPv4 N-to-N multicast tests with asynchronous traffic.

Frame size (bytes)	Maximum jitter (usec)			Throughput rate
	10% load	50% load	90% load	
64	0.060	0.090	0.720	0.110
128	0.070	0.090	0.110	0.710
256	0.690	0.710	0.140	0.710
512	0.120	0.700	0.750	0.240
1,024	0.130	0.290	0.340	0.280
1,280	0.660	0.250	0.330	0.760
1,518	0.720	0.700	0.870	1.020
9,216	0.100	0.970	1.990	1.460

Table 43: IPv4 N-to-N multicast maximum jitter with asynchronous traffic



Appendix B: Software Releases Tested

This appendix describes the software versions used on the test bed. Network Test conducted all benchmarks in September and October 2016 in a Cisco engineering lab in San Jose, California, USA.

Component	Version
Cisco NX-OS	7.0(3)I4(4)
Spirent TestCenter	4.59.7726

About Network Test

Network Test is an independent third-party test lab and engineering services consultancy. Our core competencies are performance, security, and conformance assessment of networking equipment and live networks. Our clients include equipment manufacturers, large enterprises, service providers, industry consortia, and trade publications.

Disclaimer

Network Test Inc. has made every attempt to ensure that all test procedures were conducted with the utmost precision and accuracy, but acknowledges that errors do occur. Network Test Inc. shall not be held liable for damages which may result for the use of information contained in this document. All trademarks mentioned in this document are property of their respective owners.

