

SR Redundancy and Throughput in Linux

S. Hopkins, M. Ennis
Coraid, Inc.

Summary. It is often requested that the SR be used in a configuration that permits simultaneous use of the network interfaces for aggregated throughput and redundancy. The standard solution for this is 802.3ad link aggregation, commonly known as bonding. Since AoE is its own Ethernet type and does not use IP, it has proved to be incompatible with some vendor implementations of 802.3ad bonding. To overcome this limitation, the Linux AoE driver implements a feature called multipathing to provide for round robin load balancing of multiple AoE targets on multiple client interfaces. This feature is available in 2.6 kernel Linux drivers starting with aoe6-33. Linux is currently the only system with an AoE driver implementing the multipathing feature.

There are considerations that must be made when deciding to use the second port of the SR in conjunction with multipathing. The remainder of this paper discusses issues specific to certain SR models and provides general guidelines for multipathing usage. When viewed from the rear, the first port on the SR is the left port and the second port is the right port.

SR1520 / SR420

Due to a hardware limitation in the SR1520 and SR420 series appliances, the second port provides slower throughput than the first when used for AoE. In general, Coraid recommends the second port of these appliances be used for CEC (Coraid Ethernet Console) management and not for AoE. There is no performance benefit to using both ports of these appliances with the aforementioned multipathing feature; there is in fact a performance penalty. The only benefit of using the second port for AoE is to have an additional path to the AoE storage in the event of network failure. This benefit cannot be obtained without paying the cost of reduced throughput.

The main cause for the limitation of the second interface lies in its location on the PCI bus; it must compete with the SATA cards for a limited bandwidth to memory. As a result, the second interface cannot sustain burst rates of packets as it cannot get them into memory fast enough to avoid overflowing its network receive FIFO. The result is that packets are frequently dropped. This problem is most evident when performing writes, as it is then that the receive FIFO is most stressed.

AoE drivers starting with aoe6-33 have a mechanism to reduce the number of outstanding packets a target can accept for cluster use of AoE targets. This feature helps with the receive problem of the second interface and reduces the number of retransmits that would otherwise occur.

If the SR is configured to export one Iblade and there is only one host using the SR, then the second interface limitation will only cause a minor performance penalty when used with multipathing as shown in the following tables. If, however, there are multiple Iblades being exported or there are multiple clients accessing the Iblade, then there is a potential for the second interface to heavily penalize throughput for all Iblades / users. Coraid strongly recommends not using the second interface for AoE in these configurations.

Figures A and B contain statistics achieved by averaging the results of three independent runs of `ddt` for the given SR configuration. `Ddt` is a simple tool that writes and reads to a file through a file system to determine the throughput capability of the file system and underlying storage. For a full description of `ddt`, please see the document titled SR Performance Analysis at the SR support page¹. In the figures that follow, #IF is the number of interfaces used and all throughput values

1 See Appendix A for URL

are in KiB/s. Statistics are presented using jumbo² and standard Ethernet frame sizes.

The Linux client used for these tests was stuart, a dual CPU (dual core), system with 1GiB of RAM. The Linux kernel was 2.6.18, and the AoE driver was aoe6-44. The .config file for stuart's kernel is available for download¹. The client network cards used the Intel 82546GB controller. For each configuration above, an XFS file system was placed on the resulting AoE device. The file system was mounted, and ddt was run against this mount point.

Both SR appliances were tested running the 20070111 release. For the dual port tests, two network ports on stuart were directly connected to both ports of the SR. For the single port tests, one network port on stuart was directly connected to the leftmost (first) port of the SR.

SR1521

The SR1521 appliance was introduced to overcome the limitations previously discussed with the SR1520. The SR1521 is capable of fully saturating both interfaces when used with the multipathing feature. The remaining throughput limitations are believed to be in software; Coraid is undergoing a full analysis of the SR RAID subsystem to identify bottlenecks.

Figures C and D contain ddt results for the SR1521. All configuration parameters (Linux client, SR version) are identical for the SR1521 as for the previous SR1520 tests.

² For an explanation of why we use an MTU of 4200, please see the Linux EtherDrive® FAQ entry 5.29, "Why does the Linux AoE driver only use 4132 byte Jumbo Frames?"

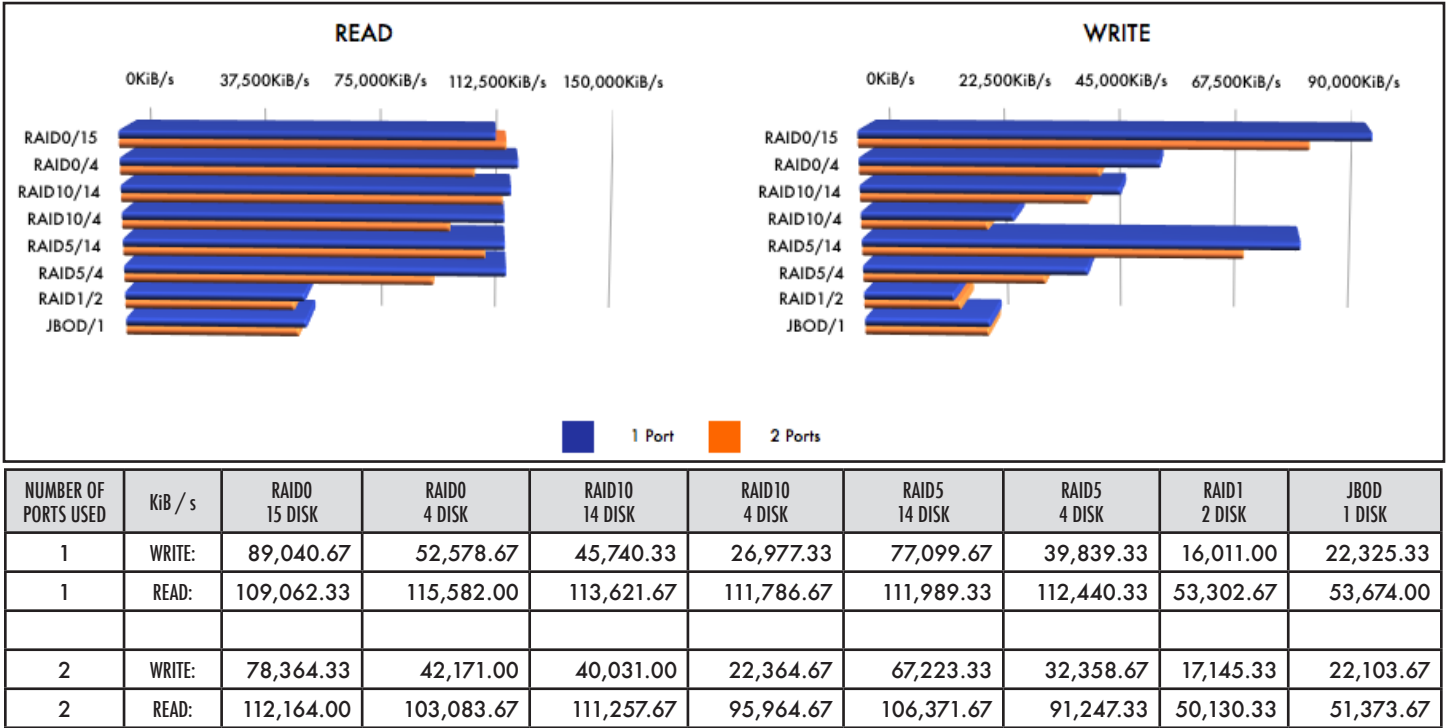


Figure A. Performance Data: SR1520, MTU of 4200 (Jumbo)

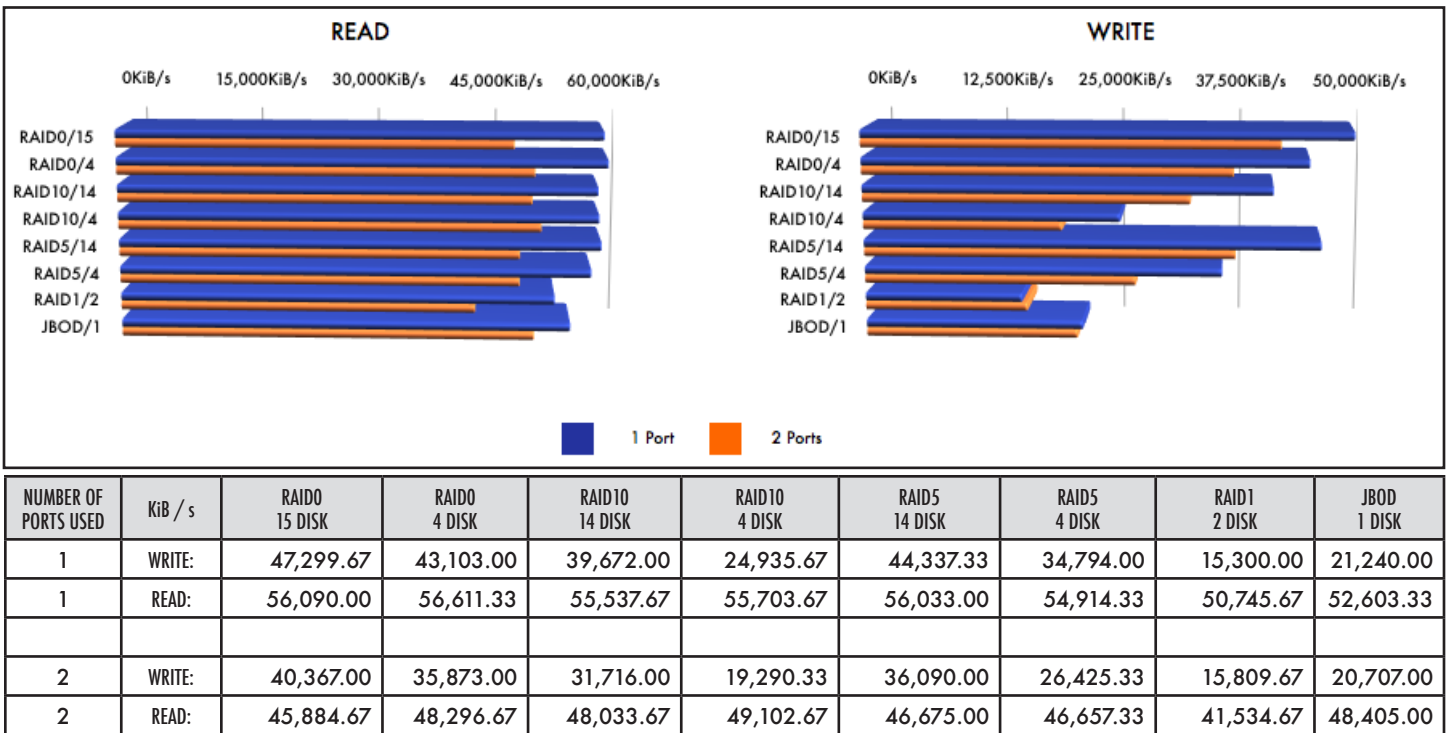


Figure B. Performance Data: SR1520, MTU 1500 (Standard)

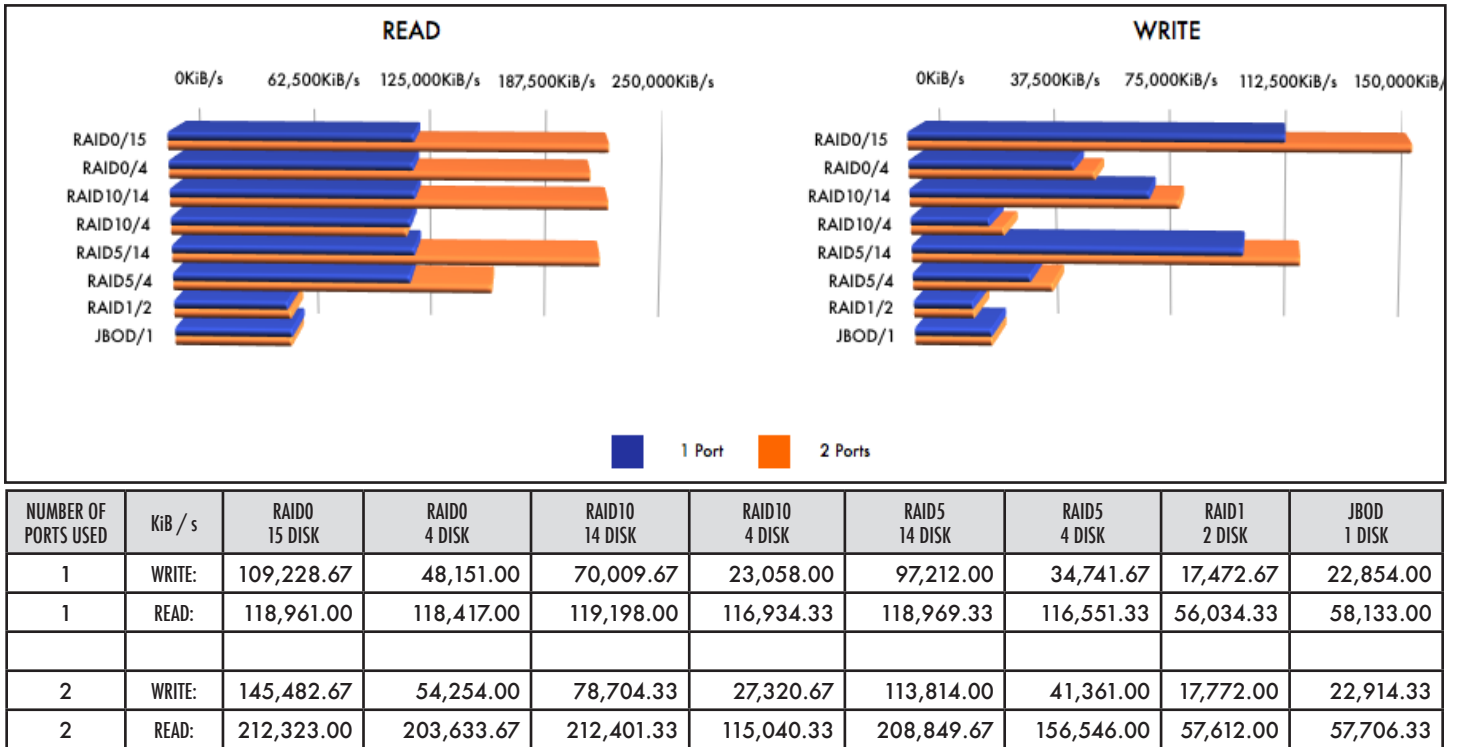


Figure C. Performance Data: SR1521, MTU of 4200 (Jumbo)

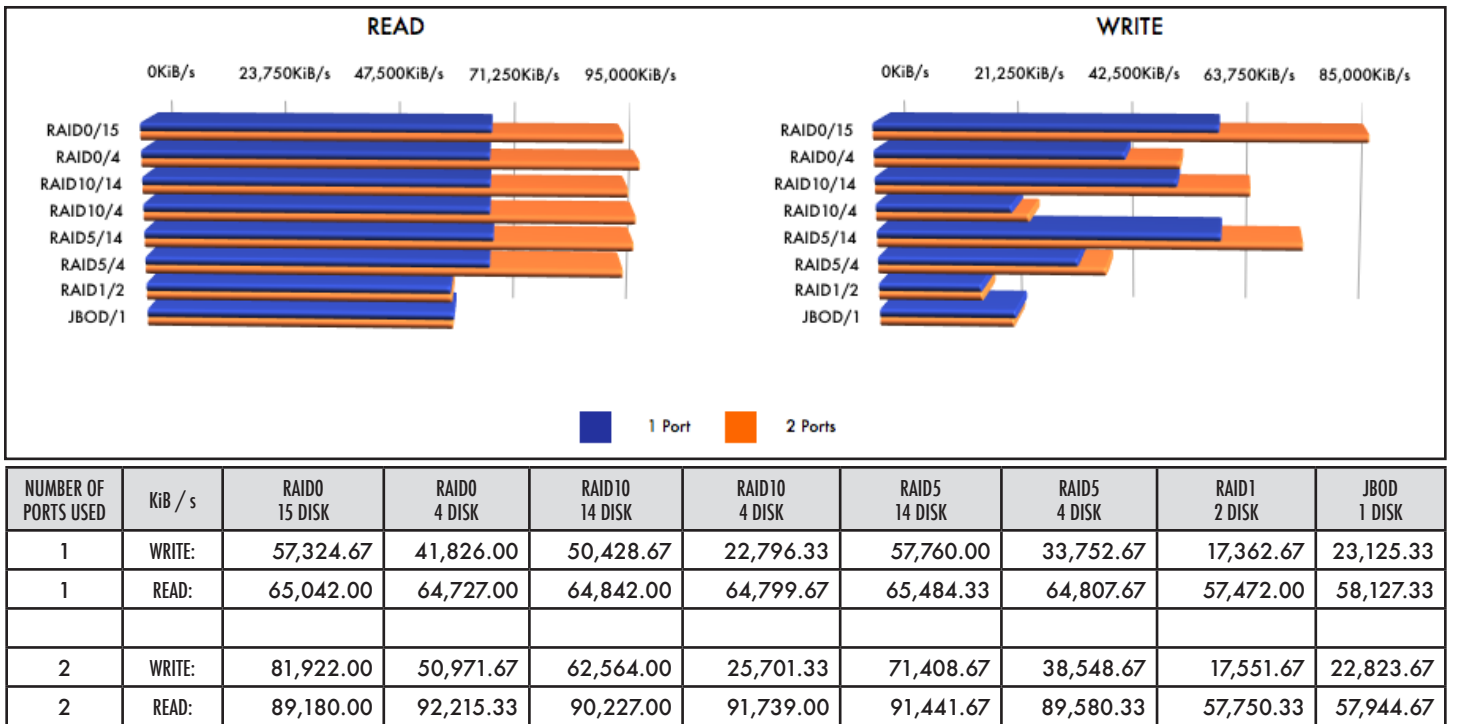


Figure D. Performance Data: SR1521, MTU 1500 (Standard)

Appendix B - References

The .config file for stuart is available at:

<http://www.coraid.com/support/sr/stuart.config>

The SR support page includes the SR firmware, user manual, and related docs:

<http://www.coraid.com/support/sr/>

Please e-mail support@coraid.com with any questions or comments.