
November 21, 2005



CASE STUDY

All Systems Down— Reprise

Edwin Hoffman, CSA

Overview

A great many people have read about the trials and tribulations suffered by the IT team at Beth Israel Deaconess Medical Center. The four days that humbled the team and the entire hospital was a severe lesson in trust, even though the team at Beth Israel was not really to blame.

For years network engineers believed that spanning tree protocol was the answer to lower cost network redundancy, and that it was a somewhat slow, but safe, system. Few knew (or believed) the horror stories that occasionally appeared that showed what a moderate sized network that is based on multiple spanning tree protect links can do when that protection starts to fail.

Excerpt from an article by Anne Barnard that appeared in the Boston Globe on November 26, 2002

Thirteen days ago, as his computer crunched the mountain of data he hoped would be his humble contribution to medical progress, the researcher - he shall remain nameless - got a phone call he'd never forget. It was Dr. John Halamka, the former emergency-room physician who runs Beth Israel Deaconess Medical Center's gigantic computer network. He told the professor that his flood of numbers was overwhelming the system, threatening to freeze thousands of electronic medical records and grind the hospital's network to a halt. "He said, 'Oh, my God!' and pulled the plug out of the wall," Halamka said last week. It was too late. Somewhere in the web of copper wires and glass fibers that connects the hospital's two campuses and satellite offices, the data was stuck in an endless loop. Halamka's technicians shut down part of the network to contain it, but that created a cascade of new problems. The entire system crashed, freezing the massive stream of information..... The problem had to do with a system called "spanning tree protocol," which finds the most efficient way to move information through the network and blocks alternate routes to prevent data from getting stuck in a loop.....

Case Statement

How does that protection start to fail? The numbers of ways are legion. Some failure points are inadvertently designed in by engineers who do not know all the rules or do not factor in traffic patterns. The main problem at Beth Israel was the spanning protection was too many levels deep. The maximum allowable is 7 levels and they had some up to 10 deep.

What does this do? It actually sets up all 10 levels to cascade into loops when a trigger event happens. What could a trigger event be? Something as simple as increased traffic flow over an open link, or a degradation of a fiber connect, or even just a bad copper cable. In Figure 1, we see a spanning tree protected network in normal operation. Bridge Protocol Data Units (BPDUs) are being issued by all three switches, but because the last topology change causes one link (blue dash line) to enter blocking mode, there is no bridge loop, and everything is fine.

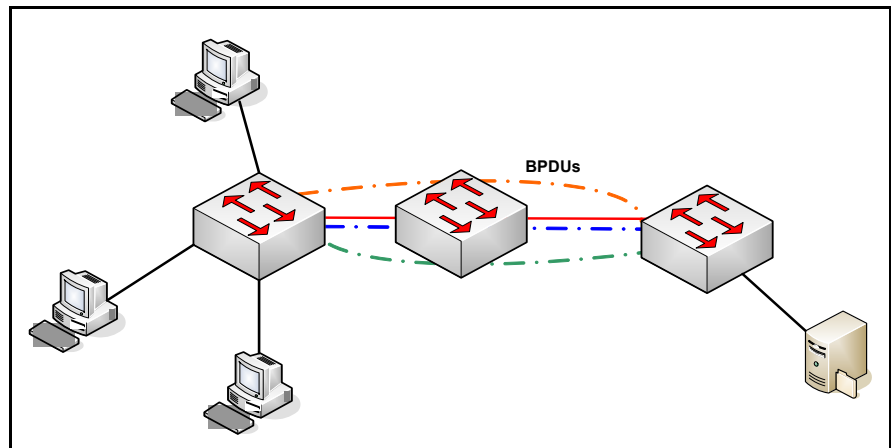


FIGURE 1. Spanning Tree Protected Network

Now traffic starts to flow and at a certain point one of the inter-switch links starts to overload and drop packets. This is a normal part of the Ethernet specification, and is completely expected, and will always happen as traffic increases. At a certain point (normally around 10–15%) BPDUs will be part of the dropped packet count, priming this network for problems.

Case Statement

Spanning Tree Protocol (STP) Failure

The primary function of the Spanning-Tree Algorithm (STA) is to cut the loops that redundant links create in bridge networks. The STP operates at Layer 2 of the Open System Interconnection (OSI) model. Bridge protocol data units (BPDUs) exchange data between bridges, and the STP delegates the ports that eventually forward or block traffic.

This protocol can fail in some specific cases, and troubleshooting the resulting situation can be very difficult, depending on the design of the network. The best way to elevate this situation is to provide routine network and maintenance troubleshooting prior to system problems develop.

A failure in the STA generally leads to a bridging loop. Most customers who require technical support for spanning tree problems suspect a system bug; but a bug is seldom the cause. Even if the software is the problem, a bridging loop developing in an STP environment originates from a port that should **block** traffic, but in a malfunction mode instead **forwards** traffic.

Figure 2 shows three workstations are working with a file server and are reaching this critical point. The increased traffic is also infecting the third switch, and it finally also submits when increased broadcast traffic (caused by the loop between the left and middle switch) reaches its port.

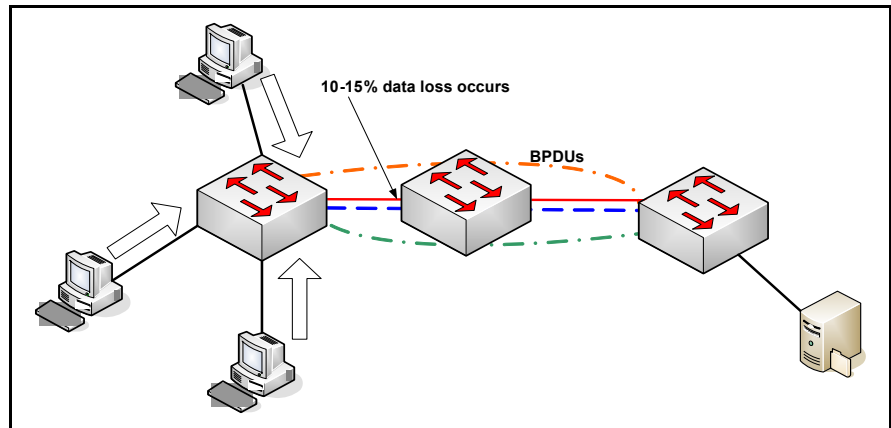


FIGURE 2. Workstations and File Server at Critical Point

Figure 3 shows the next step in this cascaded bridge loop.

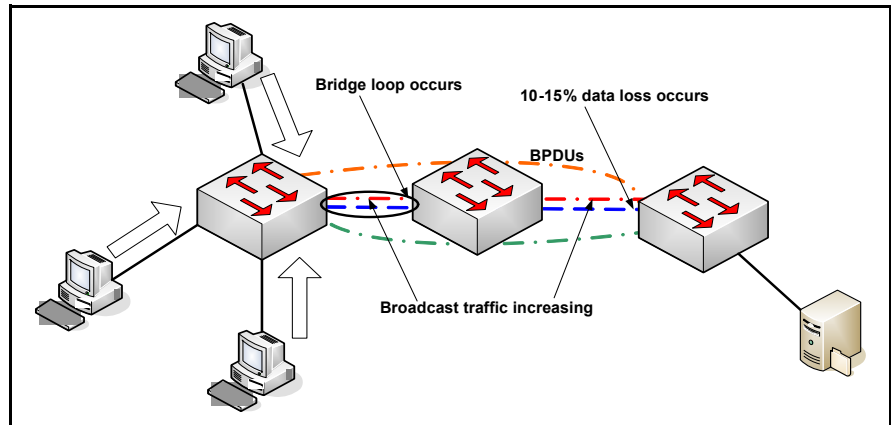


FIGURE 3. Cascaded Bridge Loop

Now the final act appears and opens the protected links between the middle and right switch, creating a bridge loop between the middle and right switch, and also creates a bridge loop between the left and right switch. The entire network is now running slow (if not completely down) and productivity is lowered.

Case Statement

Figure 4 shows this in action.

User-Visible Symptoms

How do you know when the problem has surfaced? If you're running IPX with a significant number of SAPs (services), you will completely lose connectivity with hosts on the bridged network. You won't be able to ping anything, and users won't be able to access anything. Users with slower PCs may even have problems with non-network applications like spreadsheets. IPX SAP broadcasts occur every 60 seconds, consisting of many packets, that will be issued by IPX servers as well as the routers. If you have dual-homed servers, this situation gets even worse, since the servers will send the SAP they have on network card A to network card B and vice versa. The impact of a bridging loop is magnified because the IPX networks generates much more broadcast traffic than in the normal operation mode.

In the case of IP-only networks, Raptor Networks has seen long response times, intermittent ping performance, and dropped sessions.

Source: Internet Documents

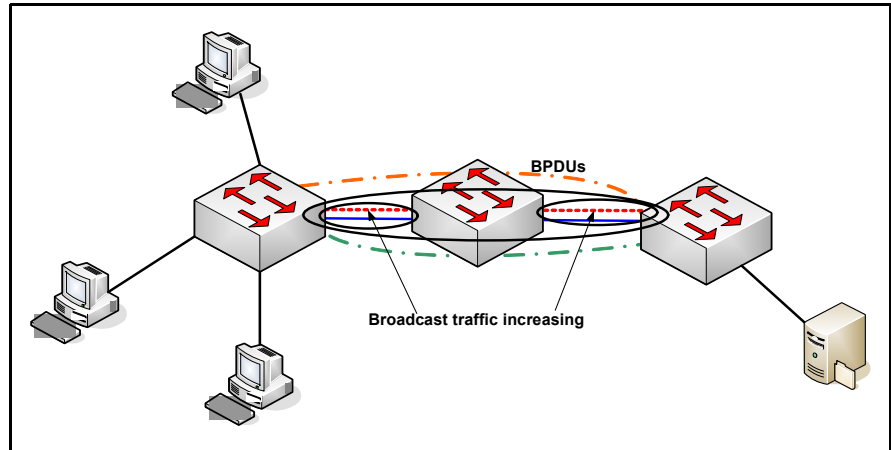


FIGURE 4. Cascaded Bridge Loop with Traffic Slowdown

So why do network engineers allow these potentially dangerous designs to be implemented?

- False trust in a technology that was designed to stop inadvertent loops being employed as a redundancy option.
- The cost of implementing a fully-switched backbone can be excessive, as evidenced at Beth Isreal, when multiple 6509s were installed to prevent this from occurring.
- A lack of knowledge of what bridge loops can do to modern networks.

Spanning tree protocol was created to prevent accidental loops in multiport bridges back in the days of when 10 Mbps was data speed, everything was half-duplex and bridges were restricted to 1-2 ports, and bridge forwarding was not wire speed. Today multiport (24, 48, 96+) switches (fast bridges) forward at a wire-speed rate of up to 10 Gbps. The potential for disaster has become greater and the potential for accidental loops has also increased. All that remains is to actually start intentionally creating bridge-loop potential and the disaster is a matter of when, not if.

How do network engineers avoid bridge loops today and still allow some redundancy? Simply with link aggregation!

The IEEE802.3ad standard has given the engineer the ability to link switches together without bridge loop potential and also increases the interswitch link speed. Several non-standard versions also exist and all prevent the need for spanning tree across multiple physical links.

The Raptor Solution

Of course, in a simple network shown in Figure 5, is the perfect way to create network- and device-level redundancy, and it will work well for most simple requirements.

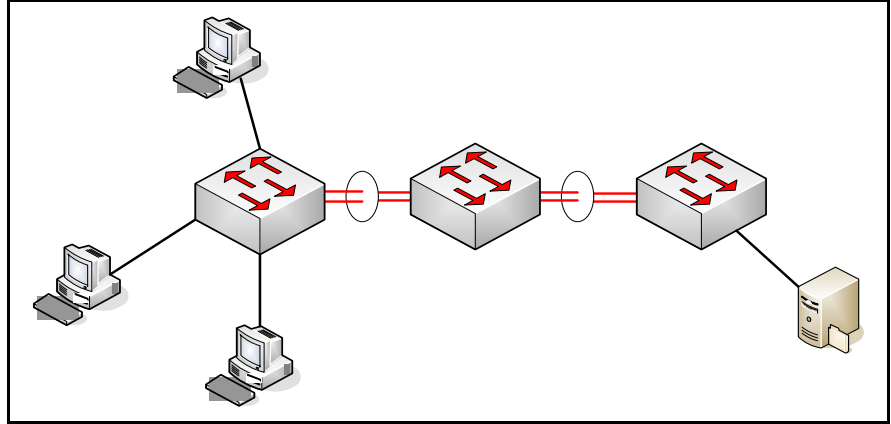
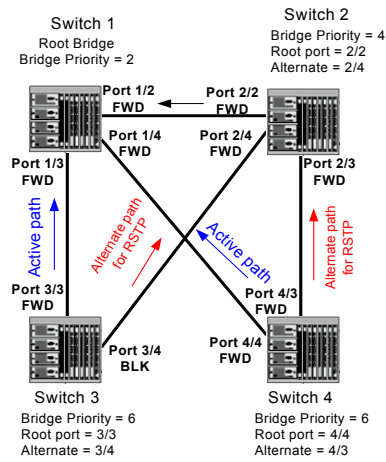


FIGURE 5. Network and Device Redundancy

RSTP Ready for Failover

A SIMPLE rapid spanning tree setup?



More complex designs bring spanning tree back into the mix. Figure 6 shows a simple change to the basic network of Figure 5 that requires spanning tree protocol to be activated.

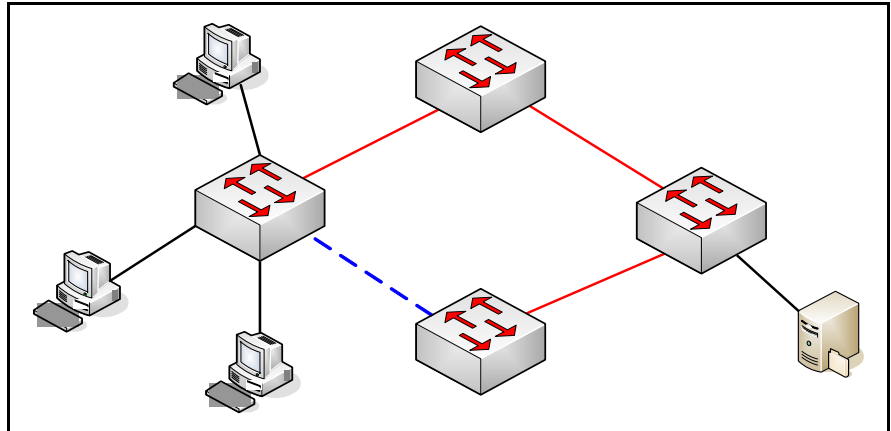


FIGURE 6. Network and Device Redundancy Requiring Spanning Tree Protocol

At this point, the specter of a bridge loop causing major network slowdowns again rears its ugly head.

The Raptor Solution

One of the main design features of Raptor Network's Ether-Raptor switch family, and in particular Raptor Adaptive Switch Technology (RAST™), was to obviate the need for using spanning tree in backbone designs.

The Raptor Solution

When the ER-1010s shown in Figure 7 are used in exactly the same way as the switches in Figure 6, spanning tree is not an issue because RAST cannot cause bridge loops to exist.

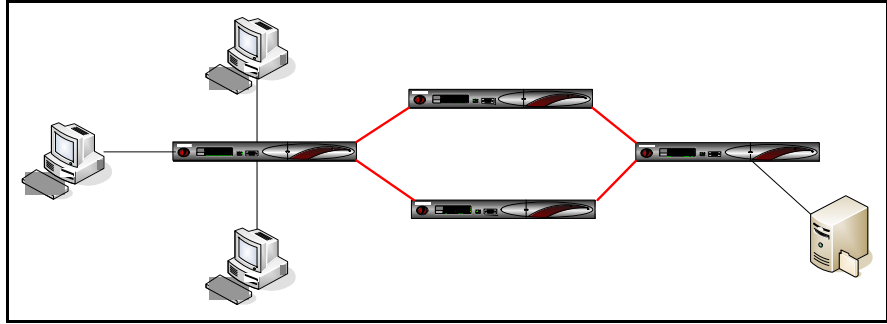


FIGURE 7. Network and Device Redundancy Without Spanning Tree Protocol

In addition, these four ER-1010 switches are actually working as a single switch unit, which allows some very useful features to be used. Now IEEE802.3ad can be used to redundantly connect other devices to this RAST-connected network in Figure 8, creating a redundant backbone that cannot cause bridge loops, but still allows levels of redundancy to exist that is beyond spanning tree's level with none of its dangerous effects. Beth Israel paid out millions to re-create their backbone—the Raptor would have cost as little at 10%, but likely not more than 25%, of the cost of their present fix.

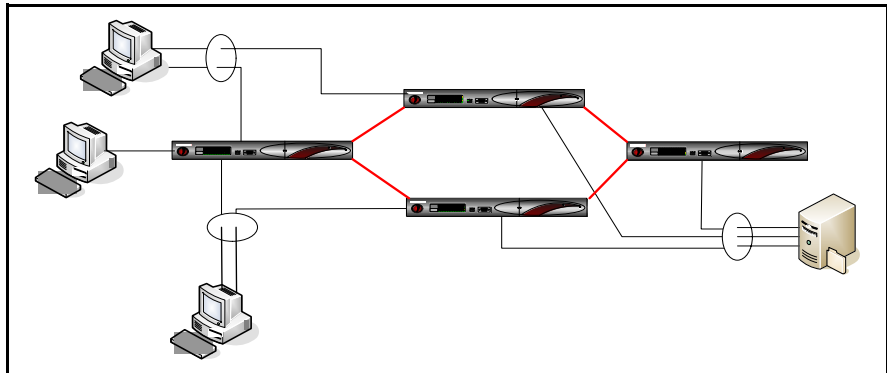


FIGURE 8. Network and Device Redundancy Using IEEE802.3ad Standard

Corporate Headquarters: 1241 E. Dyer Road, Suite 150 Santa Ana, CA 92705

Phone: 949-623-9300 / Fax: 949-623-9400 / Web: www.raptor-networks.com / E-mail: info@raptor-networks.com

Raptor Networks Technology, Inc. reserves the right to make changes without further notice to any products or data herein to improve reliability, function, or design. Information furnished by Raptor Networks Technology, Inc. is believed to be accurate and reliable. However, Raptor Networks Technology, Inc. does not assume any liability arising out of the application or use of this information, nor the application or use of any product or circuit described herein, neither does it convey any license under its patent rights nor the rights of others.

Raptor Networks Technology, Inc. is a registered trademark and RAST is a trademark of Raptor Networks Technology, Inc. All other trademarks are the property of their respective owners.