



Global Information Assurance Certification Paper

Copyright SANS Institute
Author Retains Full Rights

This paper is taken from the GIAC directory of certified professionals. Reposting is not permitted without express written permission.

Interested in learning more?

Check out the list of upcoming events offering
"Security Essentials: Network, Endpoint, and Cloud (Security 401)"
at <http://www.giac.org/registration/gsec>

The New Firewall Design Question

Jamie R. Bjerke

March 4, 2001

For IT managers, no longer is the question whether or not to implement a firewall system, the question is how to best implement a high performance, scalable, robust, highly available firewall system. Like it or not, Internet connectivity has become a mission critical resource for countless organizations. Whether providing basic Web and email services or complex e-commerce sites spanning horizontal and vertical business partners facilitating multi-million dollar transactions, the Internet has become increasingly relied upon for day-to-day business. The first firewalls were single server systems placed inline as gatekeepers for all traffic to and from the Internet. These single points of failure firewall systems generally provided the first line of defense, but at the cost of performance degradation and downtime due to hardware/software failure and maintenance. When these systems were designed and installed, performance and high-uptime were not critical requirements due largely in part to high cost/low bandwidth Internet connections and the early-on lack of dependence upon Internet services. However, as time would have it, bandwidth became less expensive and business processes more dependent upon the Internet, creating a mission critical environment for firewall systems. With this continued and increased reliance upon firewalls as the external connection point of entry for both internally and externally sourced traffic, it is with no great surprise that technologies such as hot standby, hardware-based load sharing, load balancing and clustering are becoming a prerequisite for designing a firewall system. This paper aims to define and compare these three competing technologies as well as provide example diagrams to show how each might be used to help solve the problem of high performance and/or continuous availability in firewall systems.

Hot-Standby

Hot-Standby is defined as having redundant hardware and software connected to the network in the same fashion as the production system with the ability to automatically failover all production traffic in the event of a network or system failure. Hot-Standby is generally limited to two firewalls, a primary unit which handles all production traffic under normal circumstances and a secondary or failover unit which handles all traffic when a failure occurs. Since only one unit is active at any given time, no performance enhancements are gained by implementing hot-standby.

In a hot-standby configuration the two units share Virtual IP (VIP¹) addresses as well as shared or virtual MAC addresses to handle the failover process. This technology is implemented as an add-on software component or enabled feature to firewall systems whereby intelligence is shared between the primary and failover systems with regards to the health of the network as well as the firewalls. Hot-Standby implementations typically require a heartbeat network, which handles health checking² and state³ or connection synchronization information. There are two main reasons for this:

1. The heartbeat connection can generate considerable traffic due to synchronizing state information in addition to sharing health check information. If both state and health check information cannot be successfully transmitted, the failover system can become unstable causing loss of data and or network connectivity. Case in point: the secondary unit cannot successfully communicate state or health check information with the primary, thus it deems the primary has failed. At this point, the secondary unit considers itself the primary and begins to accept network traffic. Since the primary unit has not failed, it also is still primary causing instability within the network because two systems are now answering for the same network traffic.
2. State or connection information is typically considered sensitive (since it contains IP address information as well as a picture of what traffic is allowed through the firewalls) and is better secured on a dedicated heartbeat network which is not shared by other resources.

Pros

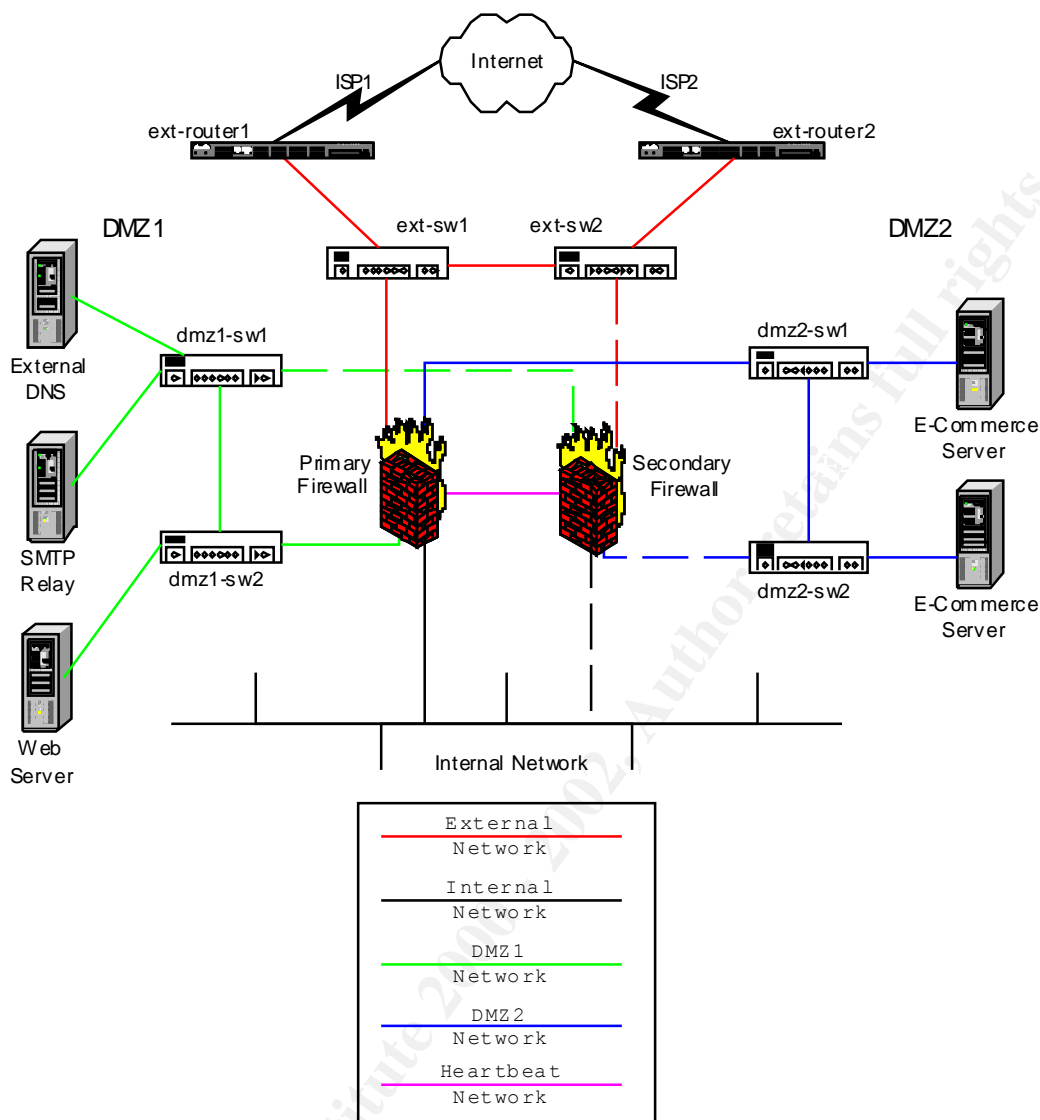
1. Generally quite easy to implement.
2. Achieves redundancy in the firewall infrastructure.
3. In addition to network health checks, certain hot-standby products support internal firewall health checks (i.e. monitoring vital processes, disk space, etc.) and will failover to the secondary unit based on robust health check criterion.
4. Maintenance and configuration changes can be made with no service disruption.
5. Most economical of the three technologies discussed.

Cons

1. No performance increases are realized. In fact, performance can degrade if the firewall system gets overburdened with processing health checks and state synchronization.
2. In many cases, the secondary or failover unit must be identical to the primary unit.
3. May or may not failover connection-oriented⁴ (otherwise known as long-lived connections) and/or Virtual Private Network (VPN) connections. These connections may need to be reestablished. This depends whether or not the firewall product used supports failover or long-lived and VPN connections.

Hot-Standby Diagram

The dashed line network connections depict the failover path through the secondary firewall should a failure occur.



Load Sharing

Load sharing can be defined as having two or more redundant firewalls, with each firewall active and passing or sharing network traffic load. This implementation does not aim to balance traffic, but does distribute the load, based on predetermined routing. A routing protocol such as Border Gateway Protocol (BGP) or Open Shortest Path First (OSPF) can be used to route traffic based on certain criteria through a set of firewalls. Routing of traffic through the firewalls will generally be symmetric⁷ as a connection originated through one firewall will remain through the same firewall unless a failure event occurs. This will create an uneven distribution of traffic through each firewall since connections will take the same path or remain static across the firewalls unless a network or firewall system event deems the traffic to be re-routed. In addition, this requires each firewall to run the routing protocol used to effectively achieve redundancy.

Generally no shared VIP or virtual MAC's are required on the firewalls to implement this design⁶.

Pros

1. Performance is enhanced since load can be distributed across multiple firewalls.
2. In addition to network health checks, internal firewall health checks (i.e. monitoring vital processes, disk space, etc.) can be custom scripted and traffic re-routing can be made based on robust health check criterion.
3. Software to perform routing protocols is inherent in routers and is freely distributed for many operating systems.
4. Maintenance and configuration changes can be made with no service disruption.

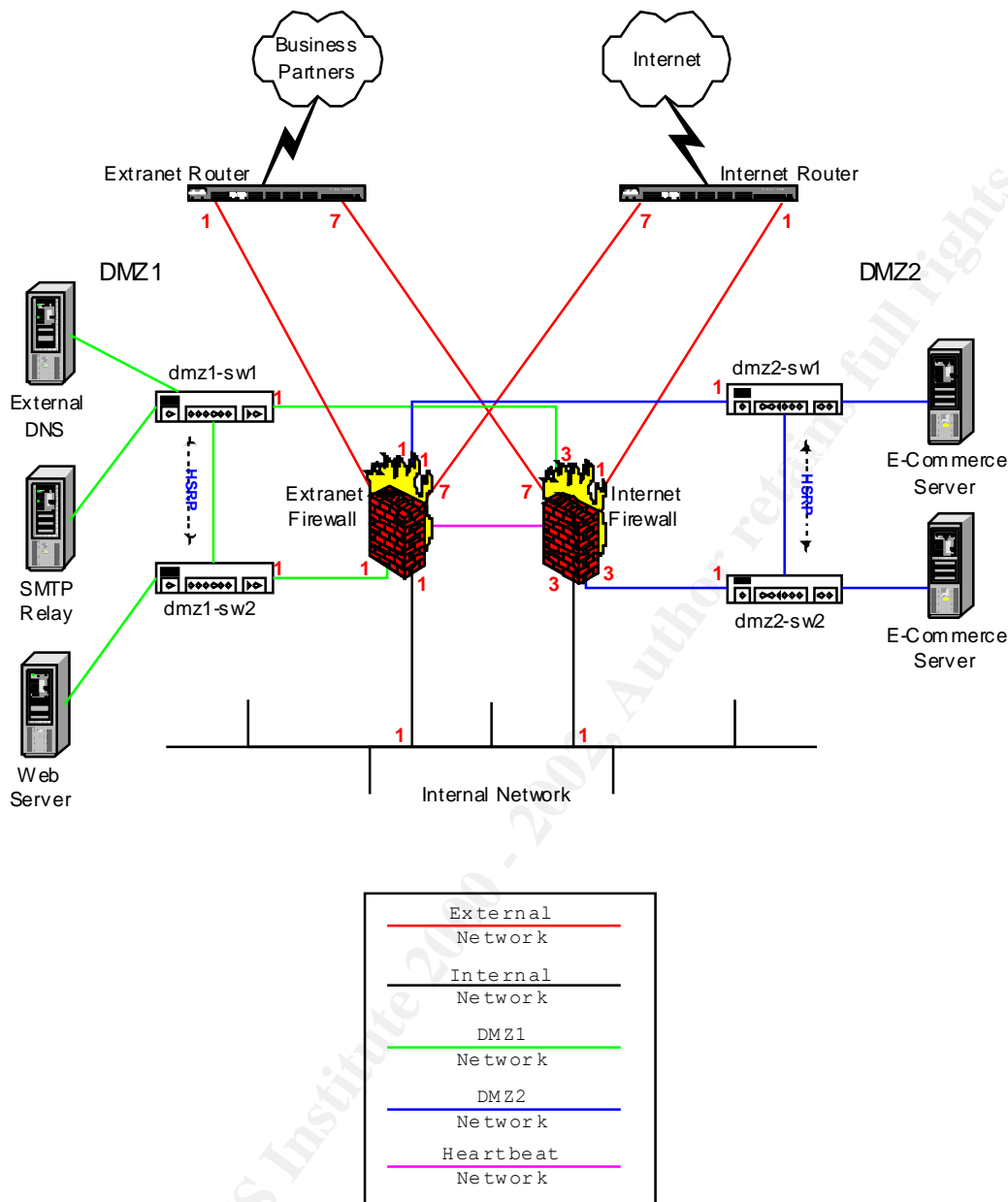
Cons

1. The load is generally not distributed evenly across the firewalls.
2. This configuration can be quite complex to implement and administer.
3. May or may not failover connection-oriented and/or VPN connections. These connections may need to be reestablished. This depends whether or not the firewall product used supports failover or long-lived and VPN connections.
4. This may not be an option when using firewall appliances since many do not support running routing protocols such as OSPF.
5. Running a routing protocol such as OSPF on the firewalls is a security risk.

Load Sharing Diagram

This diagram follows the design set forth at <http://www.hanetworks.com/networks/ospf/>
Two design differences are depicted in this diagram: utilizing two DMZ's and implementing layer three switches in the DMZ's.

© SANS Institute 2000 - 2002



Hardware-Based Load Balancing

A set of layer four switches⁵ to be used in conjunction with a set of firewall systems to provide load balancing capabilities, network health checking, network address translation (NAT) and redundancy. This configuration generally creates a near even distribution of traffic across all firewalls. However, no provisions are made to redistribute traffic under conditions where firewall load is uneven. Thus, no attempt to balance between each firewall based on current capacity benchmarks (taking into consideration free CPU cycles, free memory, etc.) are made and, therefore, only sharing of the network load results. Generally no shared VIP or virtual MAC's are required on the firewalls to implement this design.

Pros

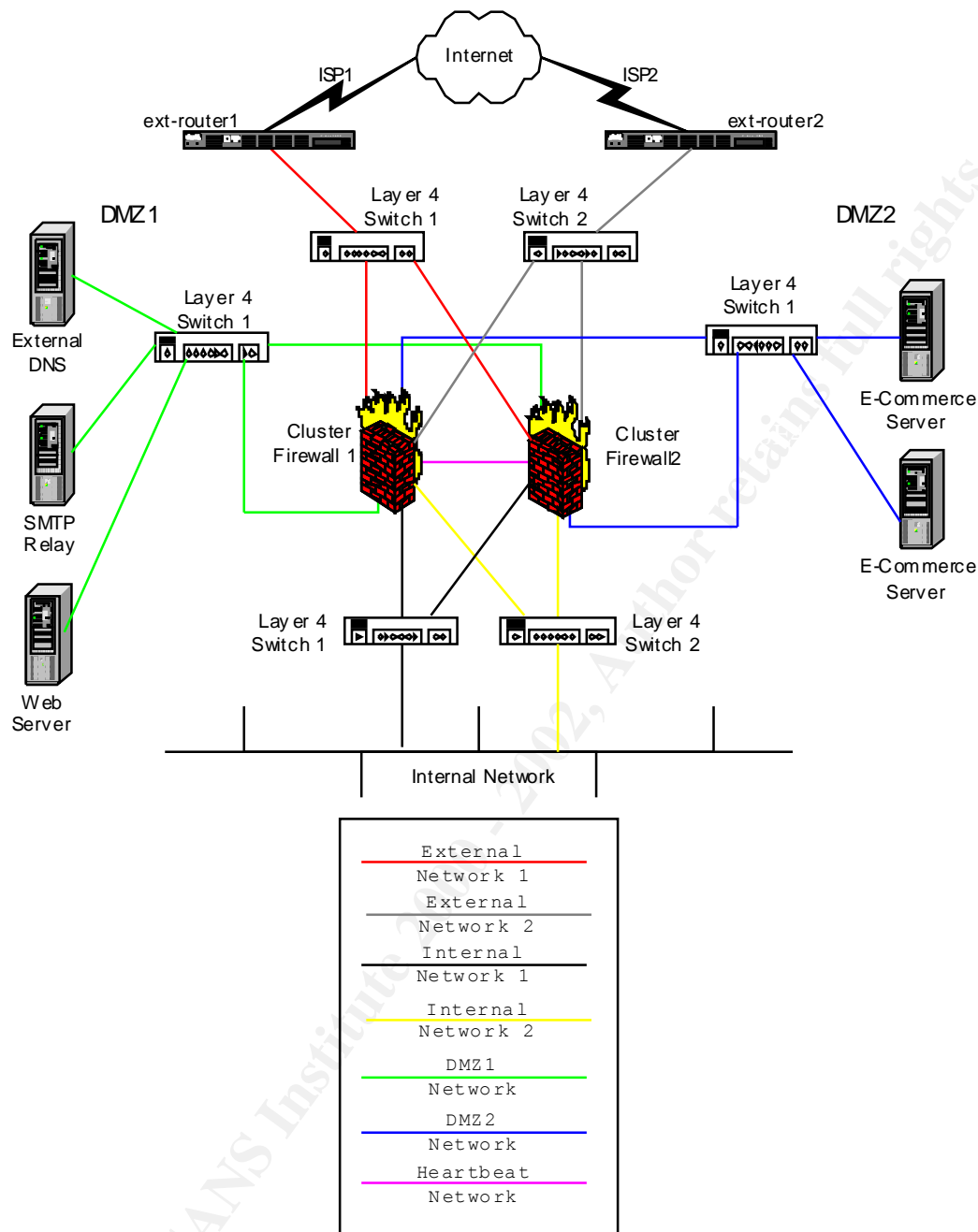
1. Performance is enhanced since load can be distributed across multiple firewalls.
2. The load is distributed near equal across all firewalls.
3. Very scalable in terms of the ease of adding firewalls.
4. Maintenance and configuration changes can be made with no service disruption.

Cons

1. A communication channel will need to be opened between the internal and external layer four switches. This means that either a rule on the firewalls will be needed to allow this connectivity or a network connection bypassing the firewalls will be required. Either is considered a security risk.
2. When running a NAT configuration, it may be necessary to have the layer four switch perform the translation. Some vendors do not support firewalls performing NAT in conjunction with layer four switching.
3. May or may not failover connection-oriented and/or VPN connections. These connections may need to be reestablished. This depends whether or not the firewall product used supports failover or long-lived and VPN connections.
4. Health checks are generally limited to network events. For instance, the layer four switch would not know if one of the firewalls was running out of disk space or memory and would continue to send traffic to this firewall.
5. Does not perform true load balancing.
6. Generally the most expensive solution when implementing full layer four redundancy in addition to firewall redundancy.

Load Balancing Diagram

© SANS Institute 2000 - 2002



Clustering

Clustering comprises all of the advantages of hot-standby and load sharing, but also adds three enhancements, intelligent/dynamic load redistribution, transparency and support for heterogeneous⁸ firewall clusters. Load redistribution allows a firewall cluster to achieve highly intelligent load balancing in addition to maintaining continuous availability. Currently, clustering is only achieved through software residing on each firewall in the cluster. Based on several metrics⁹ for each firewall within the cluster, clustering software can make intelligent decisions regarding new connections coming through the cluster as well as existing connections which may be better served by being redistributed to another

firewall within the cluster. Clustering is generally implemented using either single VIP and virtual MAC or multiple VIP's or floating¹⁰ IP address designs. Single VIP/MAC designs create transparency among the firewalls in that the rest of the network believes there is only one firewall system, greatly simplifying network designs. In contrast, the floating IP address design has the advantage of linearly increasing network throughput with each new firewall added to the cluster. Clustering software is specific to the type of firewall implemented allowing for specialized monitoring of the firewall application.

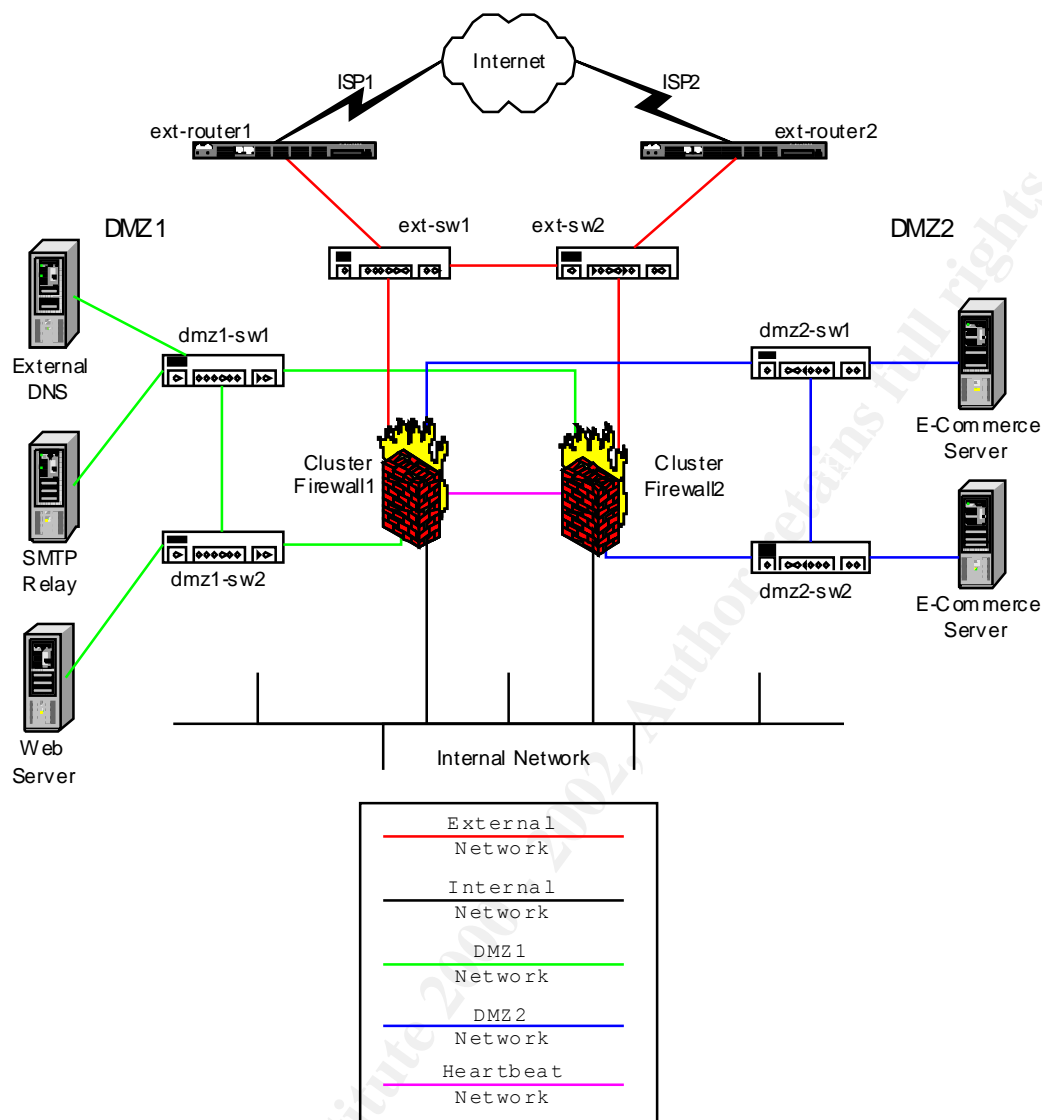
Pros

1. Unparallel performance through sophisticated load redistribution/multiple VIP capabilities.
2. The ability to define firewalls in the cluster with differing capacities (heterogeneous) to better distribute the load.
3. Robust set of health checking options including network, operating system, hardware and critical system process checks (including firewall application process monitoring).
4. A very integrated solution with all of the firewalling and load balancing components existing within the firewall servers.
5. Very scalable in terms of the ease of adding firewalls to the cluster.
6. The ability to create firewall transparency.
7. Maintenance and configuration changes can be made with no service disruption.

Cons

1. This configuration can be quite complex to implement and administer under the multiple VIP/floating IP address design.
2. Depending on the design (VIP vs. floating IP addresses) performance may be limited to the speed of the layer 2 infrastructure.¹¹
3. May or may not failover connection-oriented and/or VPN connections. These connections may need to be reestablished. This depends whether or not the firewall product used supports failover or long-lived and VPN connections.

Clustering Diagram



Sources

Cisco Systems, Inc. "Configuring Firewall Load Balancing." Content Services Switch Advanced Configuration Guide (Software Version 4.01). December 12, 2000. URL: <http://www.cisco.com/univercd/cc/td/doc/product/webscale/css/advcfggd/firewall.htm> (01 Mar. 2001).

Cisco Systems, Inc. "Advanced Configurations, Failover." Configuration Guide for the Cisco Secure PIX Firewall Version 5.3. December 20, 2000. URL: http://www.cisco.com/univercd/cc/td/doc/product/iaabu/pix/pix_v53/config/advanced.htm#xtocid57946 (01 Mar. 2001).

Foundry Networks. "Chapter 13: Configuring Firewall Load Balancing." Foundry ServerIron Installation and Configuration Guide. September, 2000. URL:

http://www.foundrynet.com/techdocs/SI/Foundry_ServerIron_Firewall_Load_Balancing.html (01 Mar. 2001).

Nokia Internet Communications Inc. "Nokia Firewall/VPN Appliances." IP Network Security Solutions. 2000. URL:
http://www.nokia.com/securitysolutions/pdf/10109_Firewall_VPN2.pdf (01 Mar. 2001).

Rainfinity. "Increasing Firewall and VPN Capacity, Performance Benefits of RainWall." January, 2000. URL:
http://www.rainfinity.com/us/eng/downloads/whitepapers/wp_increasing_fw_capacity.pdf (01 Mar. 2001).

Stonesoft Corporation. "StoneBeat White Paper." April, 2000. URL:
<http://www.stonebeat.com/dl/pdf/WhitePaper3-1.pdf> (01 Mar. 2001).

Stonesoft Corporation. "StoneBeat FullCluster White Paper." September 2000. URL:
<http://www.stonebeat.com/dl/pdf/WhitePaperFC2.pdf> (01 Mar. 2001).

Paul, Brooke, Greco, Tom, Coogan, Mike, Moore, Joel, Witty, Jason, Mesin, Alex, Angeletti, Rich. "Firewall Failover Using OSPF and HSRP." No Date. URL:
<http://www.hanetworks.com/networks/ospf/> (01 Mar. 2001).

Restuccia, Megan. "Firewall Load Balancers." November 21, 2000. URL:
http://www.sans.org/infosecFAQ/firewall/load_balancers.htm (01 Mar. 2001).

¹ Virtual IP addresses or VIP's are IP addresses shared between two or more devices. Technologies such as Hot-Standby, Load Sharing and Load Balancing commonly use VIP's as well as shared or virtual MAC addresses to "trick" the network into thinking the set of firewall systems is only one system. In a switched environment, special switch configurations are generally required to properly handle VIP's and virtual MAC's.

² Health checks range from testing network connectivity with PING to checking vital processes on the firewall. Different vendors and implementations offer varied levels of health checking.

³ State synchronization is defined as sharing TCP/IP connection information among two or more systems. For firewalls, this commonly refers to replicating the state or connection tables for TCP, UDP and VPN connections which are established through the firewall. Note: not all firewall vendors support state synchronization.

⁴ Connection-oriented connections are long-lived TCP-based connections. Examples are telnet, Citrix ICA and FTP.

⁷ Symmetric, in this scenario, means a connection going out one firewall will return back through the same firewall.

⁶ Note: there may be VIP or shared MAC addresses used on the switches and/or routers to provide redundancy among these devices. However, this is beyond the scope of the paper which only aims to discuss firewall redundancy.

⁵ Layer four switches are network devices with the intelligence to make routing decisions based on layer four of the OSI Reference Model. This means the switch can make routing decisions based on source and destination IP address as well as port numbers like 23 for telnet or 25 for SMTP.

⁸ Heterogeneous refers to the ability to have different hardware among the firewalls within the cluster. For instance, one firewall could have two CPU's and 1 GIB of RAM while the other has only 1 CPU and 256MB of RAM. The main requirements are that the systems have the same architecture and operating systems. While heterogeneous configurations like the one mentioned above are technically possible, they

are not recommended configurations by the product vendors. The real world implementation of this would be defining different capacities for firewalls with similar, but not equal processing power.

⁹ Metrics are measurements of certain operating system and/or application events such as CPU cycles, memory usage, disk usage and buffer utilization.

¹⁰ Floating IP addresses are implemented by assigning a number of virtual IP addresses equal to the number of firewalls in the cluster. These IP addresses are then moved around from firewall to firewall according to the load balancing software to balance the traffic load.

¹¹ A detailed discussion of the layer two limitations when using a single VIP is beyond the scope of this paper. However, this information can be found at

http://www.rainfinity.com/us/eng/downloads/whitepapers/wp_increasing_fw_capacity.pdf

© SANS Institute 2000 - 2002, Author retains full rights.