

Load Balancing for Microsoft® Office Communication Server 2007® Release 2

A Dell™ and F5 Networks® Technical White Paper

End-to-End Solutions Team
Dell | Product Group – Enterprise

Dell/F5 Partner Team
F5 Networks



THIS WHITE PAPER IS FOR INFORMATIONAL PURPOSES ONLY, AND MAY CONTAIN TYPOGRAPHICAL ERRORS AND TECHNICAL INACCURACIES. THE CONTENT IS PROVIDED AS IS, WITHOUT EXPRESS OR IMPLIED WARRANTIES OF ANY KIND.

© 2009 Dell Inc. All rights reserved. Reproduction in any manner whatsoever without the express written permission of Dell, Inc. is strictly forbidden. For more information, contact Dell.

Dell, the DELL logo, PowerEdge, PowerVault, and Dell EqualLogic are trademarks of Dell Inc. *Microsoft* is either a trademark or registered trademark of Microsoft Corporation in the United States and/or other countries. *EMC* is the registered trademark of EMC Corporation. Other trademarks and trade names may be used in this document to refer to either the entities claiming the marks and names or their products. Dell Inc. disclaims any proprietary interest in trademarks and trade names other than its own.

CONTENTS

INTRODUCTION	4
MICROSOFT OFFICE COMMUNICATIONS SERVER 2007 R2 (OCS R2)	5
LOAD BALANCING FOR OCS R2 COMPONENTS.....	5
HARDWARE LOAD BALANCERS.....	6
F5 BIG-IP® LOCAL TRAFFIC MANAGER (LTM).....	7
FEATURES	7
OCS R2 LOAD BALANCING BEST PRACTICES	10
HIGH AVAILABILITY (HA)	10
SSL ACCELERATION	11
TCP IDLE TIMEOUT	11
LOAD BALANCING AND PERSISTENCE	11
SNAT	12
CONCLUSION	13
REFERENCES	13

Introduction

Microsoft® Office Communications Server 2007 Release 2 (OCS R2) allows enhanced instant messaging, presence indicators, conferencing, email, voice mail, fax, and voice and video communication to take place from a Communicator and Outlook end point through the data network. Higher workloads on enterprise unified communications deployments that are associated with these features may require the use of multiple servers that perform the same function. In order to distribute TCP traffic to the OCS R2 servers evenly in such enterprises, load balancers should be deployed.

To support larger deployments of Microsoft Unified Communications and for production purposes, it is recommended to setup a hardware load-balancer for scalability and high availability of the Edge server, Front End server pool, Director and Communicator Web Access server. For the scalability of other server roles, Microsoft Network Load Balancer (NLB) can be used, if required.

This whitepaper describes the OCS R2 infrastructure components and how they can be effectively scaled for enterprises. As an example, F5 Networks BIG-IP® Local Traffic Manager (LTM) in an OCS R2 environment is presented to illustrate best practices for configuration and hardware load balancing.

Microsoft Office Communications Server 2007 R2 (OCS R2)

Rich voice and video communication can take place from a computer, OCS enabled phone, smart phone, or pocket-pc using Microsoft Office Communications Server 2007 RS (OCS R2). OCS R2 can also be connected to the public switched telephone network (PSTN) from a mediation server attached to an IP-PBX or VOIP gateway. This allows calls to be forwarded to VoIP phones, traditional telephones, or cell-phones connected to the PSTN allowing “anywhere” access to OCS users. In OCS R2, new features such as group chat, team calling, response groups, dial-in conferencing, and SIP trunking are introduced.

Most of these OCS features have unique requirements, and often deployed on separate servers to provide higher quality of service and user experience. Depending on the load and configuration requirements, some components are deployed on multiple servers to balance the network traffic and provide a greater degree of availability. These configurations may require different load-balancing techniques to properly distribute the load.

Load Balancing for OCS R2 Components

Load balancing can be achieved using software, hardware or DNS round-robin methodologies. Hardware load-balancing is more connection oriented and is the most robust load-balancing mechanism. Unlike software or DNS round-robin technique, hardware load-balancing is not native or integrated as part of the application. Four components of the OCS deployment do require hardware load-balancers for larger and scalable configuration. These components include Front-End servers (Enterprise Edition), Communicator Web Access servers, Directors, and Edge servers. Clients that logon from the internal network register with the Front-End servers. Front-End servers also handle instant messaging and routing of calls. There are also a number of other services collocated with the Front-End servers, including A/V, IM, telephony, and web conferencing. The standard edition Front-End server does not support load balancing and is recommended only for smaller deployments.

The Communicator Web Access (CWA) server allows a user to access Office Communicator features using a supported Web-browser over a secure channel. Edge servers that are load balanced will direct Session Initiation Protocol (SIP) traffic to the Director array. A Director within the array will then authenticate the requests made from the external network, and also route the traffic to the Front-End server pool. Other components of OCS R2 can be made highly-available using the built-in load balancing mechanism. Figure 1 displays a load-balanced topology.

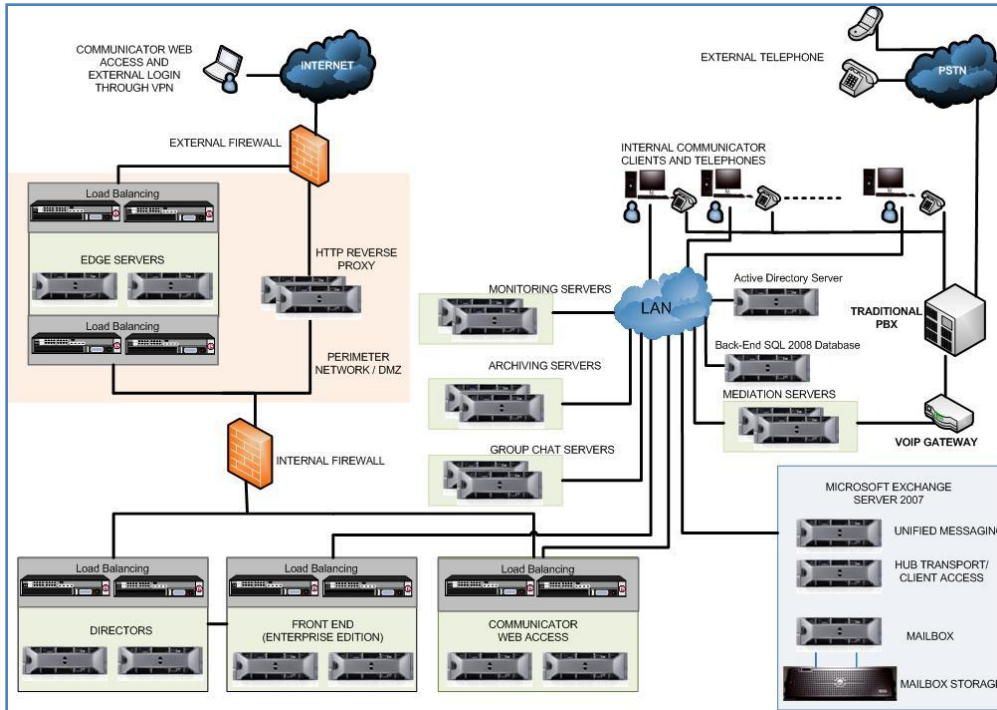


Figure 1: MS-OCS R2 Architecture with Load Balancing

The internal network is protected from Internet attacks by an internal and external firewall. This network houses all the essential services necessary for unified communications. The internal network also shows the Director, Front-End and Communicator Web-Access components within the hardware load-balancer topology. The DMZ contains the Edge servers that in turn run three services called the Access Edge service, A/V Edge service, and Web Conferencing Edge service. The external network allows Communicator Web Access, Remote User Access, Federation and Public IM connectivity.

Hardware Load Balancers

Hardware Load balancers distribute application traffic, such as SIP and HTTPS, across a number of servers. They increase the capacity and fault tolerance of applications and improve overall application performance by decreasing the server overhead associated with managing application connections, traffic optimization, and encryption off-load. Requests that are received by a load balancer are distributed to a particular server based on a configured load-balancing method. Some industry standard load-balancing methods are round robin, weighted round robin, least connections and least response time. Load balancers ensure reliability and availability by monitoring the "health" of applications, and only sending requests to servers and applications that can respond in a timely manner. As shown in Figure 1, there are some OCS R2 roles that are not presently supported for hardware load balancing. For these roles, built-in

mechanism of load balancing can be considered. The advantages of load balancers are the following:

- Even distribution of Communicator and Communicator Web Access traffic to the OCS R2 infrastructure.
- Uneven loads can be distributed across OCS R2 servers that have different processor and memory capabilities.
- Scalability for rapidly growing unified communications deployments.
- High availability for enterprises; each load balancer serves multiple OCS R2 servers and distributes the traffic between them. In addition, hardware load balancers can be made highly available with sub-second fail-over capabilities so that they are not single point of failure.
- Protection against various attacks, such as denial of service (DoS) and distributed denial of service (DDoS), from the Internet on OCS infrastructure. This functionality is the most useful for the external-facing load balancers deployed with the Edge server array.
- SSL offload and acceleration. CPU cycles on backend servers, such as CWA, are conserved by ensuring that encryption and decryption of Web-traffic is done by load-balancing components.
- Security. Communicator and Communicator Web Access clients cannot directly connect to back-end OCS R2 infrastructure. In addition, the hardware load balancers are default-deny devices; only services that are explicitly allowed in the load balancer's configuration would be filtered through its firewall.

F5 Big-IP® Local Traffic Manager (LTM)

F5 BIG-IP® Local Traffic Manager™ (LTM) turns your network into an agile infrastructure for application delivery. It's a full proxy between users and application servers, creating a layer of abstraction to secure, optimize, and load balance application traffic. This gives you the control to add servers easily, eliminate downtime, improve application performance, and meet your security requirements.

LTM makes in-depth application decisions without introducing bottlenecks for the Communicator clients, and ensures that only explicitly allowed services can pass through to the OCS R2 application servers. The following sections explain LTM features and best practices for building a reliable Application Delivery Network.

Features

A variety of LTM features are necessary to load balance OCS R2. The features of LTM can be accessed using a secure Web browser connection (HTTPS) to the IP address of the management port. The Web interface provides access to the system configurations, application configurations and monitoring. Figure 2 displays the LTM interface.

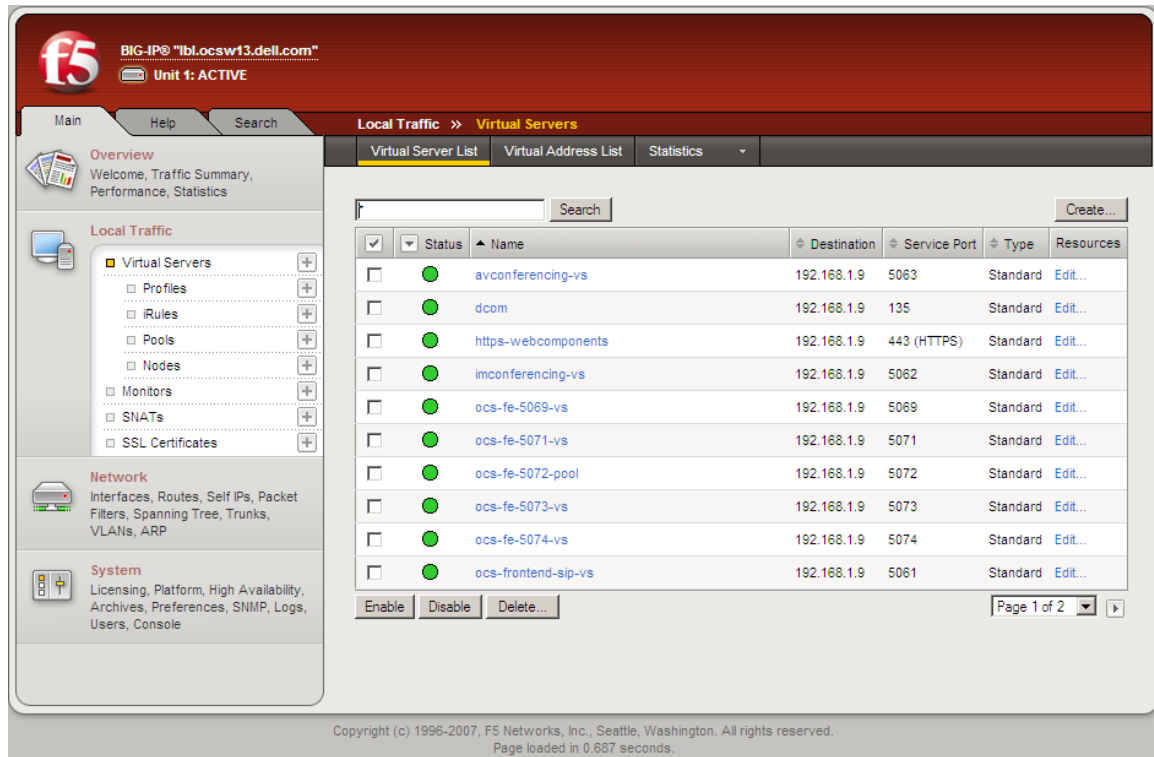


Figure 2: LTM secure Web interface

Some of the key features of the load balancer are virtual servers, profiles, SNAT and pools. The virtual server is an IP:Port combination (VIP) for the load-balanced application or service with an assigned pool. The DNS name represents the OCS R2 pool FQDN and should be associated with the VIP. A different virtual server will need to be configured for each application service that is published. For example, a virtual server supporting front-end SIP communications running on 192.168.1.9:5061 with a DNS name of pool.ocsw13.dell.com. For a different front-end service, such as the response group service, a new virtual server will be set up for 192.168.1.9:5071 using the same DNS name. SNAT and profiles, such as persistence or SSL acceleration, can be assigned to the virtual servers. The pool definitions include the creation of pool members, health monitors, and selection of the load-balancing method. For example, CWA pool member 192.168.1.18:443 has a HTTPS monitor and a load-balancing method of “Least Connections (node)”. Table 1 provides load balancing feature definitions for OCS R2.

Table 1: Definitions of load balancing features

Virtual Server	IP:Port combination representing a load balanced application or service. An object managing Virtual IP Address (VIP), Service Port, Profiles, Resource Pool, Connection Mirroring, SNAT, iRules and other traffic management features.
Profiles	Objects managing Protocol and Persistence settings related to network traffic.
Pool	Object managing Load Balancing Methods, Health Monitors and Pool Members.
Health Monitors	Objects managing health check parameters and monitoring processes for application services.
Load Balancing Methods	Objects managing the load distribution to the application servers. Stop sending traffic if a member is marked down or fails to respond to health checks.
Pool Members	IP:Port combination representing the services running on one or more servers, virtual servers or other devices.

Figure 3 displays two Communicator clients sending requests to the Front-End server pool. The first client is load balanced to Front End 1 server. The second client logs in shortly afterwards and is load balanced to Front End 2 server using the Least Connections (node) algorithm. Other clients that log in later will be distributed across the Front-End servers according to the load balancing method.

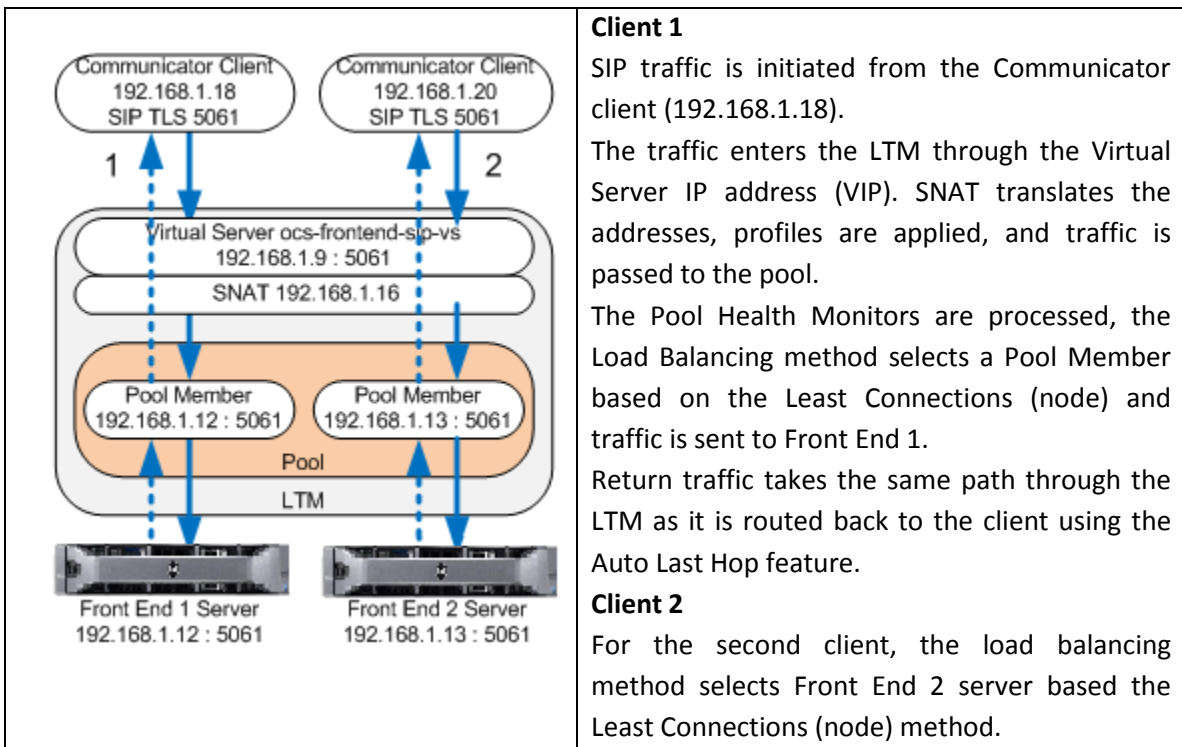


Figure 3: Communicator client load balancing

For OCS R2 hardware load balancing to work correctly, certain ports and protocols must be allowed to be load balanced on the LTM. Table 2 lists the ports that need to be assigned to the virtual servers.

Table 2: Ports and Protocols Used by the Load Balancer for Front-End Servers
(referenced from the Microsoft Web Site)

OCS R2 Load Balanced Component	Port	Protocol
Front End Servers	5060/5061	TCP MTLs
Front End Servers	443	HTTPS
Front End Servers	444	HTTPS
Front End Servers	135	DCOM and RPC
Front End Servers	5065	TCP
Front End Servers	5069	TCP
Front End Servers	5071	TCP
Front End Servers	5072	TCP
Front End Servers	5073	TCP
Front End Servers	5074	TCP
Communicator Web Access server	443	HTTPS
Director	5060/5061	TCP
Edge Servers	443	TCP
Edge Servers	5061	TCP
Edge Servers	5062	TCP
Edge Servers	3478	UDP
Edge Servers	443	TCP
Edge Servers	5061	TCP
Edge Servers	3478	TCP

The Front End, Director and Edge servers must be enabled for TCP or MTLs traffic through ports 5060/5061. The virtual server is essential for signaling using Session Initiation Protocol (SIP). The CWA virtual server uses port 443 for load balancing Web traffic from the Communicator Web Access clients.

OCS R2 Load Balancing Best Practices

The following sections describe the OCS R2 best practices for LTM load balancing.

High Availability (HA)

High Availability (HA) mode provides application fault tolerance for OCS R2. HA requires two LTM units and additional configuration to complete the set up. These configurations can include failover cabling, MAC masquerading, configuration synchronization, and connection mirroring. Benefits include LTM redundancy and sub-second failover in the event of a network outage, power outage, or failure on the active LTM. This means that Communicator, Communicator Web Access and external users of OCS R2 receive uninterrupted service should one of the LTMs fail.

In the case of an active LTM failure, the standby LTM detects the failure of the active LTM and initiates a failover. LTM #2 becomes the active unit and resumes all Communicator requests, as shown in Figure 4.

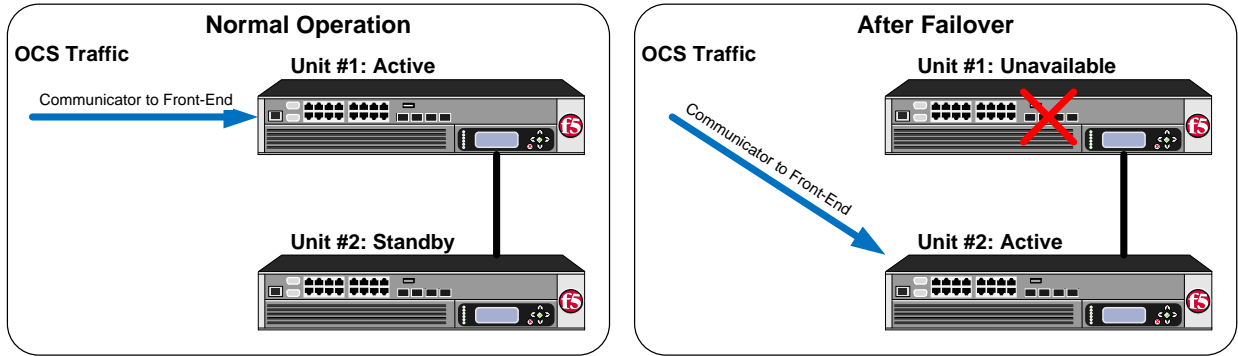


Figure 4: LTM High Availability (H/A), Failover

SSL Acceleration

SSL acceleration can be implemented on the CWA servers to off-load the HTTPS encryption processing and reduce CPU utilization. In order to enable SSL acceleration for Communicator Web Access servers, first enable the CWA servers to process HTTP (instead HTTPS) and then load a certificate on the LTM. Associate an SSL profile with the virtual server, and when the CWA clients connect to the CWA virtual server the LTM will decrypt the incoming HTTPS requests and load balance them to one of the CWA servers in the pool.

TCP Idle Timeout

This value specifies the number of seconds that a TCP connection to a virtual server can remain inactive before the load balancer will close the connection. For Communicator Web Access servers, an idle timeout of 1800 seconds is recommended. For Front End server load balancing, this value should be set to 1200 seconds.

Load Balancing and Persistence

The persistence feature ensures that once a client is load balanced to a pool member, all future connections from that client are directed to the same member. The Least Connections (node) load balancing method and HTTP cookie insert persistence are recommended for CWA virtual servers. This allows session affinity between the Communicator Web Access clients and the CWA servers. Simple persistence can be used to ensure that SIP over TLS connections to a Front-End server in the Enterprise Pool maintain affinity.

SNAT

The SNAT feature can map multiple Communicator or CWA client IP addresses to a translation address defined on the LTM. The translation address can be a single SNAT address, a SNAT Auto Map address, or a SNAT pool. A SNAT pool should be used when handling more than 65,000 simultaneous connections. When the LTM receives a request from a client IP address, and if the client IP address in the request is defined in a SNAT, the LTM system translates the source IP address of the incoming packet to the SNAT address. SNAT is a requirement for load balancing OCS R2 Enterprise Front-End server pools.

Conclusion

Office Communications Server R2 can provide rich voice, video, and live meeting capabilities for the enterprise. For large deployments, keeping these services highly available and provide application scalability requires the use of hardware load balancers for the Front-End server pool, CWA servers, Edge servers, and Directors. The F5 BIG-IP® Local Traffic Manager™ (LTM) is an easily-configurable hardware load balancer with an HTTPS administrative interface. As a best practice for deploying OCS R2, the LTM should be made highly available by using an active/standby configuration and SSL acceleration be enabled for CWA secure Web communications. When deployed to provide high availability for Director and Edge server roles, external access through the Internet can be made fault tolerant. Persistence and idle timeout should also be enabled on the LTM to ensure that the connections are properly managed.

References

Dell Unified Communications:

<http://www.dell.com/unified/>

F5 Networks Deployment Guide for Microsoft OCS:

<http://www.f5.com/pdf/deployment-guides/microsoft-ocs-ltm94-dg.pdf>

F5 Networks Products Overview:

<http://www.f5.com/products/>

Microsoft Office Communications Server 2007 R2 Deployment Guide:

[http://technet.microsoft.com/en-us/library/dd441174\(office.13\).aspx](http://technet.microsoft.com/en-us/library/dd441174(office.13).aspx)