# Active-Active and High Availability

Advanced Design and Setup Guide

Perceptive Content Version: 7.0.x

**perceptive**software

from Lexmark

# Table of Contents

# Active-active server configuration and a highly available system

This document explores an active-active server configuration and a highly available system as they relate to Perceptive Content Server and Perceptive Content Client components. Both of these systems enable you to configure multiple servers so that you increase your system's reliability and keep downtime to a minimum. Previous versions supported active-passive failover. Now you can implement an active-active system for continual client access when a failover occurs. You can also set up a high availability cluster to provide continued service when system components require service or fail.

This document is divided into the following sections:

- Active-active configuration for Perceptive Content Server

- High availability in an Perceptive Content system

This document is designed to provide an overview of these configurations, as well as an architectural model, to help understand how to implement these configurations in a Perceptive Content system. This document addresses an audience familiar with server architecture or server administration.

For more specific details about failover and active-active server configuration for Perceptive Content, refer to the following documentation list.

| Reference Document | More Information |
|---|---|
| *Perceptive Content Server and Client Installation and Setup Guide*<br><br>If you are updating your Perceptive Content system from a previous version, also refer to the *Perceptive Content Server and Client Update Readme*. | Details information about installation and updating existing versions. |
| *Perceptive Content Technical Specifications* | Provides information about requirements of the system. |
| *Perceptive Content Architecture Technical Paper* | Provides overall architecture information. |
| *Perceptive Content Active-Passive Failover Cluster Administrator Guide* | Describes active-passive failover. |
| *Perceptive Content Server Cluster Resource Monitor Installation Guide* | Provides additional information about the Microsoft Cluster Server as part of your high-availability system. |

# Active-active configuration for Perceptive Content Server

In Perceptive Content 7.0, you can install and configure Perceptive Content Server as active-active, with embedded services installed on multiple single-server instances or nodes. The database, OSM, and shared files reside as separate services.

In an active-active environment, you configure and run two or more active servers. These servers run at the same time as individual servers with separate workloads. In the event that one server fails, another server picks up the workload and client connections from the failed server without interrupting transactions or affecting end users. As soon as the failed server is back online, it reclaims the transactions from the secondary server and resumes its duties.

**Note**  Before setting up your system for high availability, refer to the documentation for any third-party hardware and software products you are using.

## Benefits of an active-active server setup

To create system redundancy in an active-passive server configuration, an environment is installed on one server and additional environments are installed on companion servers. These servers run on separate nodes, sharing the workload between them. This companion server sits unused most of the time, only being utilized when the primary server fails. This server configuration requires that a user launch a new session and log in again after services have been restored. Agents must also reconnect and establish sessions when the primary server becomes unavailable.

When you setup an active-active server configuration for Perceptive Content 7.0, two or more servers run simultaneously. When a server becomes unresponsive, the system automatically allocates resources to one or more available server instances. In this scenario, clients automatically reconnect to the new server, while resubmitting uncommitted transactions. Users running client applications remain unaware of the server switch. This active-active server switch occurs seamlessly and without the need for you to start a second server at the time of failure.

Load balancing is an inherit feature of an active-active server configuration. As a result, you can distribute workloads using load-balancing hardware devices, which allow the distribution of the workload across multiple server machines. Using load balancing, you can quickly increase efficiencies in your system performance.

An active-active server setup also allows for easier system maintenance, because the system continues to handle transactions when an administrator takes a single server offline for routine maintenance.

## Running multiple instances of Perceptive Content Server

The following sections contain information about running multiple instances of Perceptive Content Server, including distinguishing log and temporary file outputs, IP addresses, and node paths. The following sections use the following terms: real server and server farm. In this definition, a real server is a dedicated physical server that you can group to comprise a server farm.

### Server-instance naming

During the installation of an instance of Perceptive Content Server, the instance provides a unique name to identify that instance within the system. Additionally, instance names are unique across all instances of the Perceptive Content Server. The instance names distinguish log and temporary file outputs in environments that have multiple instances installed to a clustered file system. Configuration files utilize instance names so that multiple instances can share the configuration files. Instance names are often required when performing certain server management activities.

## Traffic routing

Traffic routing for an active-active Perceptive Content Server setup is available for passing multiple server IP addresses to the server. For example, an agent attempts to connect to each IP address in order, until the agent finds an open address to connect to the server.

## Using nodes

In an active-active environment, nodes actively route data just as they do in a traditional server environment. However, in an active-active environment, there is one node instance per machine, spread across several machines. When you configure nodes, ensure that paths in the database are valid on both nodes and that settings in INI files are valid on both nodes.

## Using a server farm for load balancing

A server farm is a collection of physical servers that operate behind a virtual IP address, streamlining the server workload by spreading it among many physical servers using a load-balancer. For example, when a connection is made to a virtual IP address with a load-balancer, the load-balancer picks the best real server to handle the connection. A server farm also increases redundancy by allowing other servers to handle incoming requests if one fails.

The following is a high-level overview of setting up a server farm. For detailed instructions, refer to the "Assemble and configure a server farm" section in the *Perceptive Content Server and Client Installation and Setup Guide*.

# Server health monitoring

Health monitoring is the system a load-balancer uses to determine if a physical server is available to accept incoming connections. It includes a simple machine ping to determine if the machine is online, and a specific Perceptive Content Server probe that verifies that the server is running and responding to requests. Use the following parameters to set up health monitoring.

**Note**  The following overview, which is detailed in the "Set up server health monitoring" section of the *Perceptive Content Server and Client Installation and Setup Guide* contains steps that are specific to the Cisco Application Control Engine (ACE) Module. When setting up your system's health monitoring, refer to the third-party documentation for configuring steps specific for the third-party product you use.

- Determine the type of health probe. For Perceptive Content Server, the type is a TCP health probe.

- Set the probe interval count, which is the time interval between sending probes during a health check.

- Set the pass-detect interval, which is the time interval between sending probes during a health check when the server is in a known bad state.

- Set the pass-detect count, which is the number of successful responses a probe must produce before the server is marked as healthy.

- Set the fail-detect count, which is the consecutive number of times a probe must fail before the server is marked as failed.

- Set the timeout-response count, which is the amount of time that a server has to return a response during a probe. If the server does not return a response within the set time, the probe fails and the server is marked as failed.

## Physical servers

A physical server is a machine that hosts data, manages network resources, and processes the workload from clients. Adding physical servers involves setting up the servers on a VLAN interface. The servers are located on the same subnet as the VLAN interface. After setting up servers on machines, use the following steps to add them to the hardware load-balancer configuration:

- Name the servers.

- Set the IP address.

- Set the connection limits.

- Start Perceptive Content Server on the real server machines.

## Configure a server farm

Configuring a server farm involves initiating the health probe and setting an action for a failed health check. You can load balance Message Queuing Agent. For more information on configure Message Queuing Agent, refer to *Perceptive Content Server and Client Installation and Setup Guide* for your database.

Use the following parameters to set up the server farm:

- Set purge as the action for a failed health check.

- Add the real server to the server farm.

- Set the port number. This is the port where Perceptive Content Server listens, and is the same port that is set for the health probe.

- Set a backup server and port. This server becomes active if the real server is in a failed state.

- Set the server weight. Servers with higher weights receive more connections as a ratio of their weight to other servers' weights.

- Deploy the server.

## Configure a virtual server

Virtual servers are interfaces that accept incoming connections and route them to a real server. Use the following parameters when setting up a virtual server:

- Set the virtual server IP address. This is the address the client uses to connect to one of the machines in the server farm.

- Set the VLAN interface.

- Set load balancing as the primary action.

## Command line arguments

You can use service commands or INTool to configure active-active in the system using command-line arguments. The following sections contain a high-level overview of service commands and INTool. To learn more about using command-line arguments to manage the system in an active-active configuration, refer to "Administer the Server" in Administrator Help.

### Use service commands

You use service commands to manage services, such as Perceptive Content Server and agents. With service commands, you can display the status of a server instance, start and stop specific instances of a server, and install or uninstall server instances. You can also use service commands to manage multiple instances of Perceptive Content Server.

The service command syntax is `<service> <-command> <instance_name>`, where `<instance_name>` allows you to execute a command on multiple instances of the same server that run in parallel in an active-active environment.

### Use INTool

INTool is a command-line tool that is installed by default when you install Perceptive Content Server. Using INTool, you can manage your server, databases, and object storage manager data. You can also rate your system performance, manage licenses, and obtain table structure information.

To display a list of INTool commands, open a command prompt window, change the drive to the inserver\bin directory for 32-bit system, or inserver\bin64 for 64-bit system, and then type `intool`.

## High availability in an Perceptive Content system

**High availability is** a component of an active-active server environment that, through hardware and software configurations, enables a process or application to continue operating by eliminating single points of failure and minimizing disruption to end users. The following information provides a high-level overview of a Perceptive Content Server system with high availability.
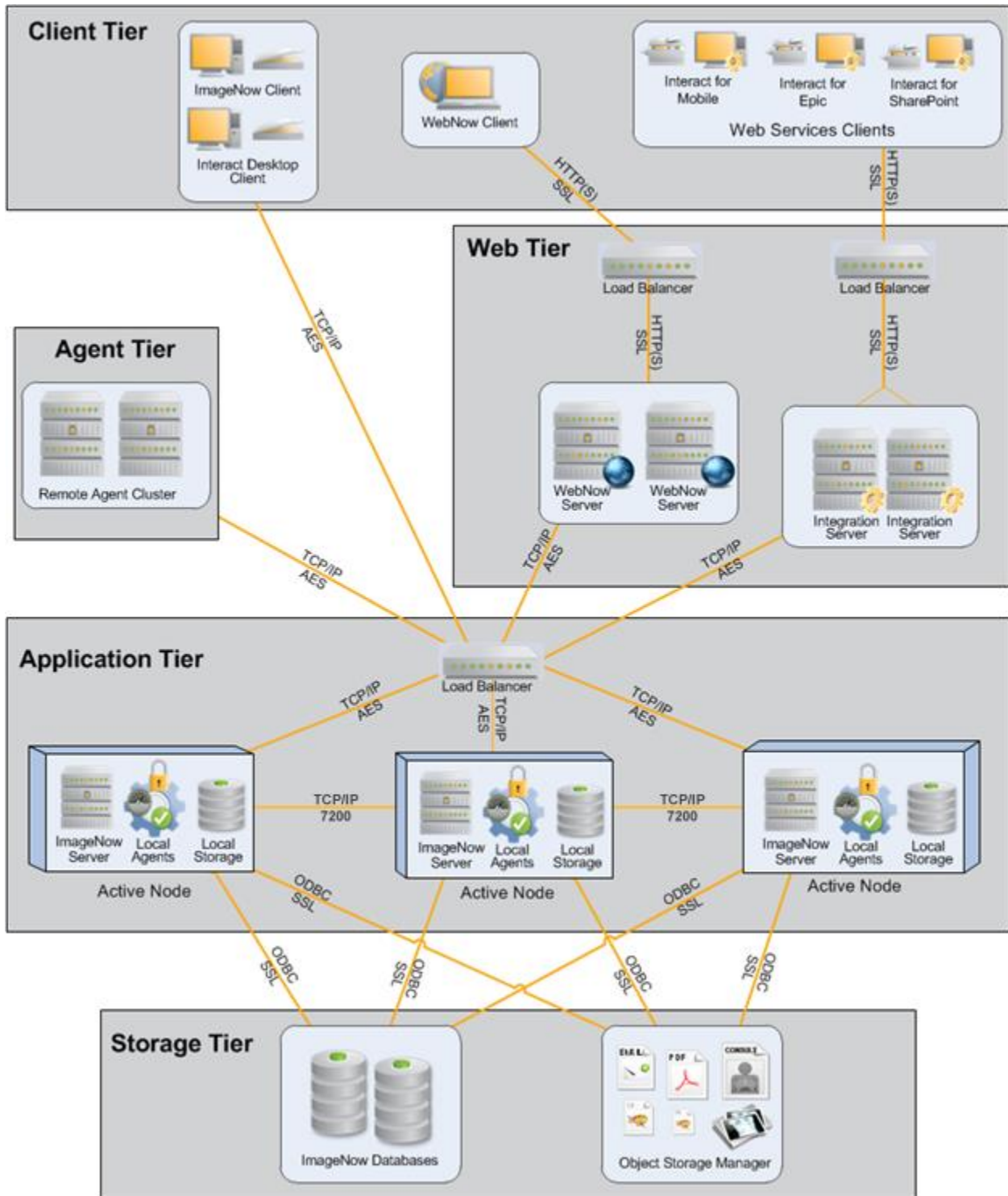
You can configure Perceptive Content servers in a load-balancing cluster configuration to achieve high availability. In a load-balancing cluster, multiple servers run simultaneously while a load balancer—a third-party device or software—balances calls to the system between available servers. A clustered configuration includes at least two Perceptive Content Server nodes, an active database cluster, shared SAN or NAS storage containing the OSM, database files, and server directories.

Perceptive Content Server also operates proactively when configured for high availability. When you balance multiple instances of Perceptive Content across redundant nodes, health policies ensure that you can detect and address problems with a server before they become severe. For example, when a health policy finds an issue on a server node that could disrupt access to an application, the instructions in the policy can redirect it so that the node no longer participates in the system and direct its traffic to other nodes. Therefore, you do not need to take the system down to repair it.

## Perceptive Content high availability system architecture

The following table and figure is a high-level overview of the system architecture behind a high-availability Perceptive Content Server configuration.
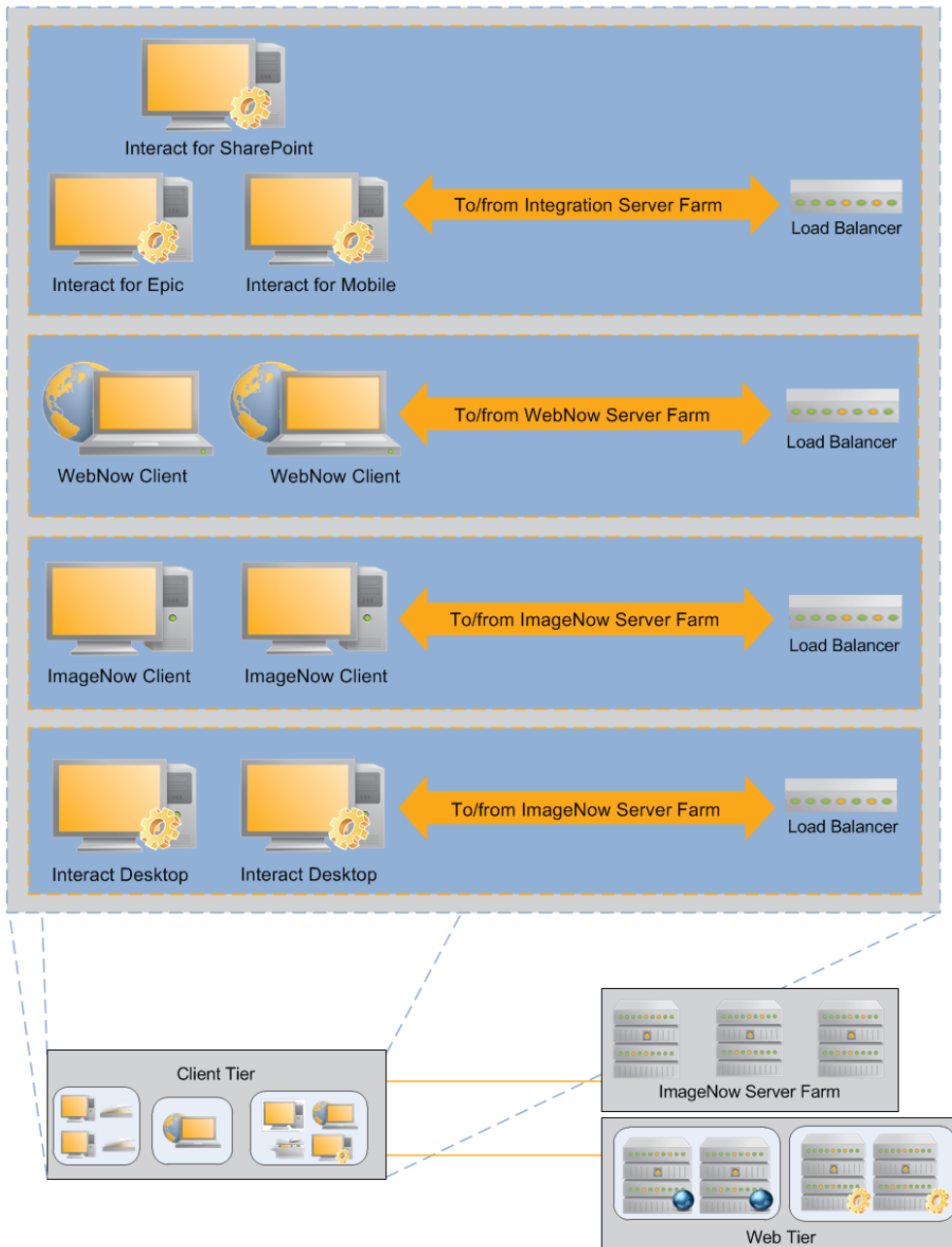
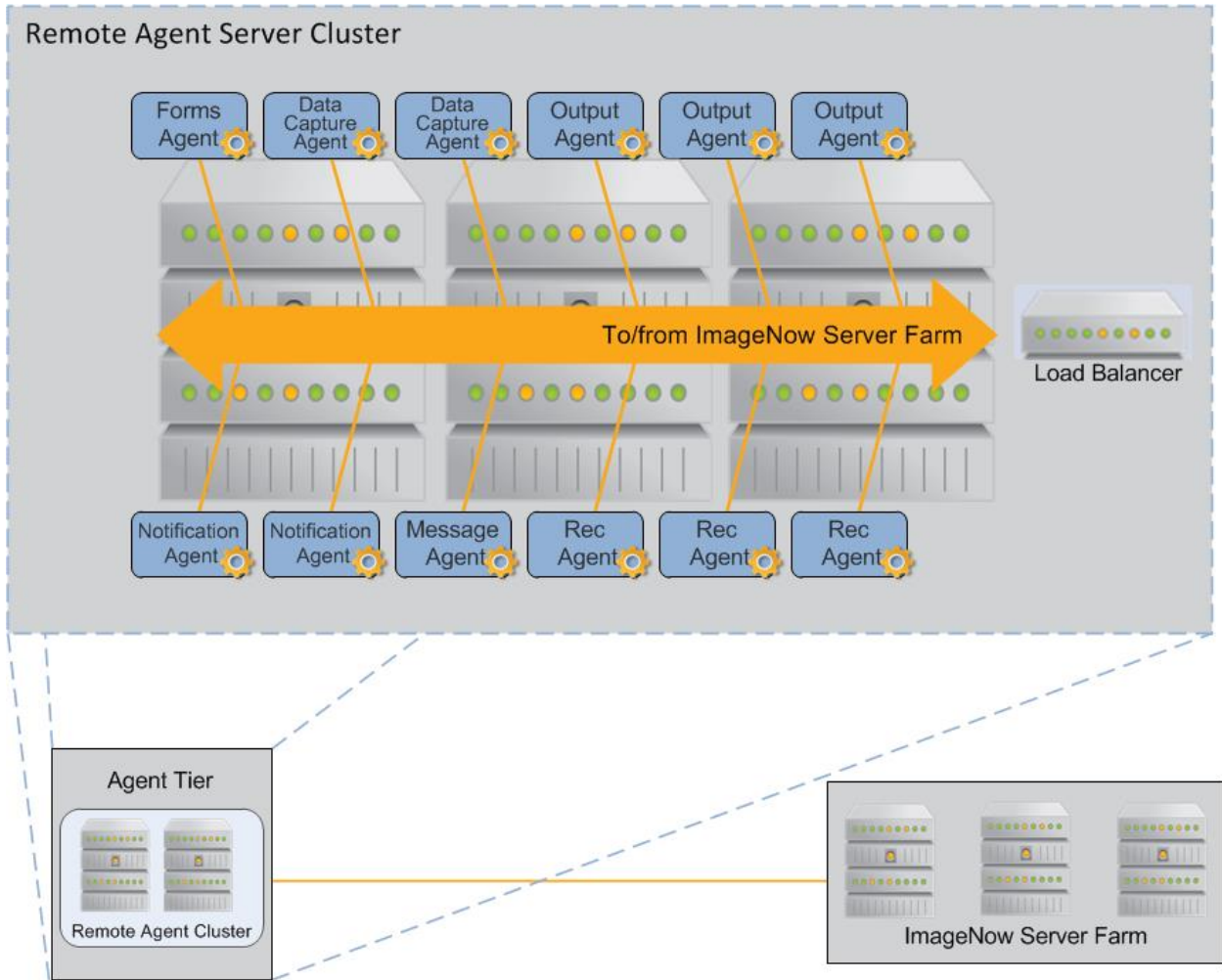| Tier | Description |
|------|-------------|
| Client | Contains the Perceptive Content Client, Interact Desktop, WebNow Client, and other web clients such as Interact for Epic and SharePoint. |
| Agent | Contains external, server-side agents, such as Forms Agent, Output Agent, Mail Agent, and Recognition Agent. Additional instances of these agents might provide enhanced performance along with optimizing your overall system. Refer to the agent installation guide for more information. |
| Web | Contains the products that connect Perceptive Content to the web, including WebNow Servers and Perceptive Integration Server. |
| | WebNow accesses data from the Perceptive Content Server through a web application server. Integration Server allows for third-party application functionality that is compatible with HTTP web services, to send and receive data from Perceptive Content Server. |
| Application | Contains one or more nodes, where each node contains an instance of Perceptive Content Server. Also contains local storage for the bin, log, and temp directories and embedded agent directories, such as job and workflow. |
| | **Note**  In an active-active configuration, the nodes must directly connect via port 7200. The Message Queueing Agents do not communicate through the load balancer. You may notice interconnection errors and some unresponsive services if the nodes cannot access port 7200. |
| Storage | Contains the Perceptive Content database, which stores the metadata and system information, and the Object Storage Manager (OSM), which stores the document objects. |

## Client Tier

The Client Tier contains Perceptive Content Clients, WebNow clients, Interact Desktop, and other Web service clients such as Interact for Epic and SharePoint. In this tier, communication occurs as follows:

- For Perceptive Content Client and Interact Desktop, communication occurs using TCP/IP and AES directly from a load-balancing device in the Application tier.

- For the WebNow clients, communication occurs using HTTP or HTTPS and SSL directly from the load-balancing device in the Web tier. HTTP or HTTPS and SSL manage communications between the web service clients and the Web tier.
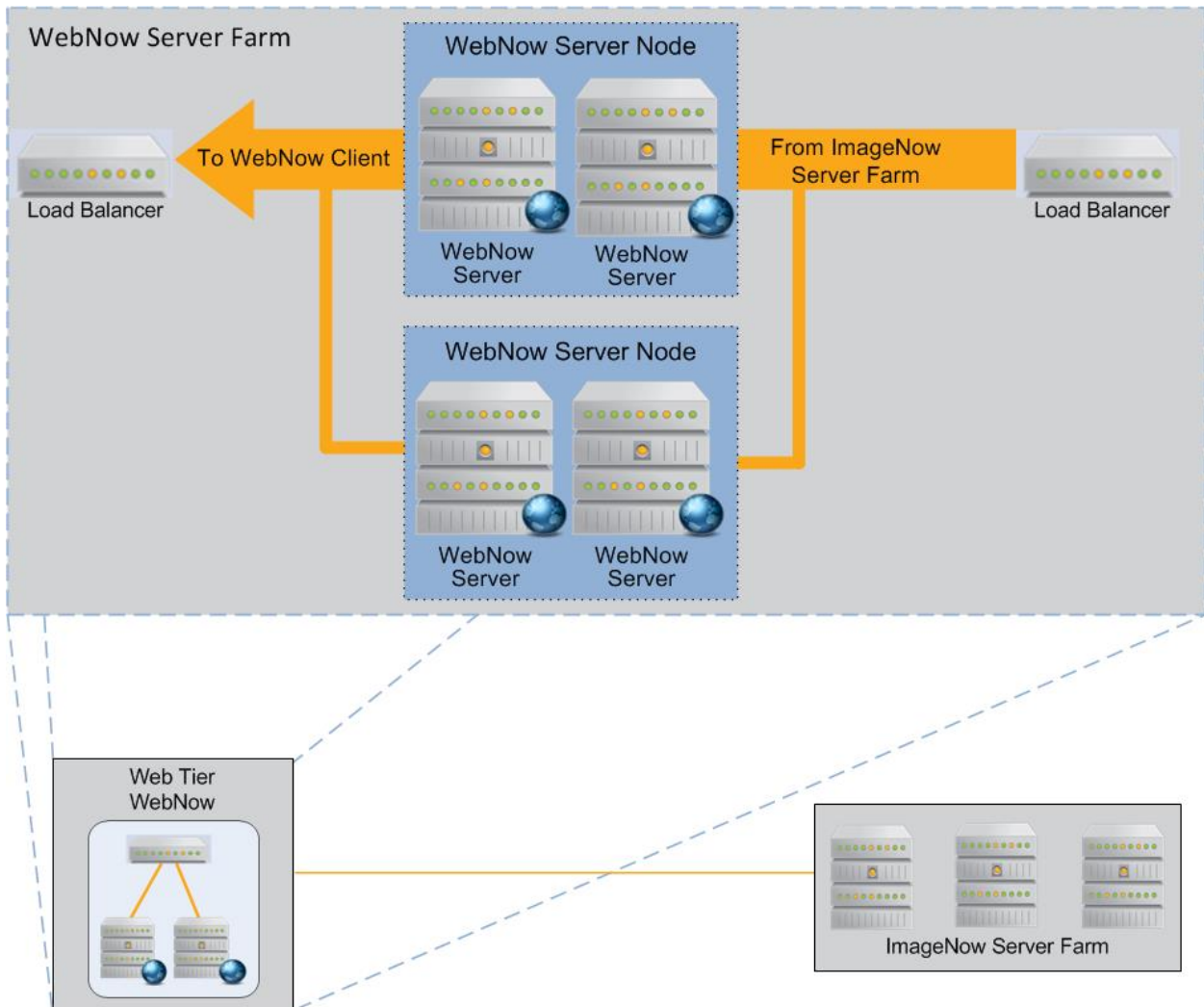
## Agent Tier

The Agent Tier contains any agents you have running on remote agent servers, like the agents listed in the top row of the following figure. Remote agents can run in an active-active, clustered configuration and connect to a Perceptive Content Server running on any supported operating system via TCP/IP. Communication between remote agents and the Client Tier routes through the load-balancing devices in the Applications Tier and the Web Tier.
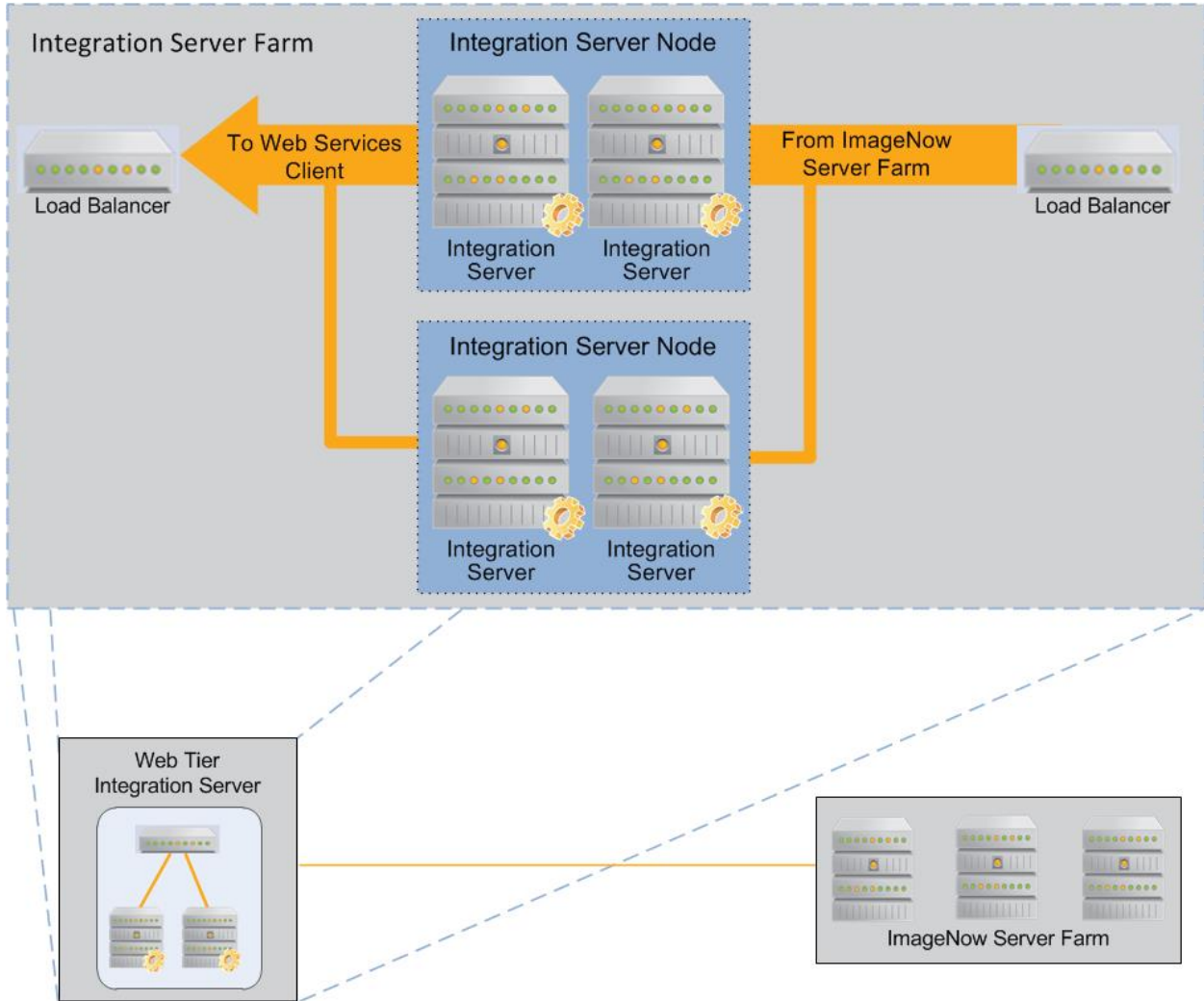
# Web Tier

The Web Tier contains the products that connect to Perceptive Content through the web, including WebNow Servers, which enables communication between Perceptive Content Server and WebNow client machines. The tier also contains Integration Servers that enable seamless communication between Perceptive Content Server and third party applications that are compatible with web services.

You can combine more than one WebNow Server to form a cluster, and combine two or more clusters to form a farm. Incoming communication occurs through TCP/IP and AES and outgoing communication occurs using HTTP or HTTPS and SSL. Incoming and outgoing traffic routes through load-balancing devices.
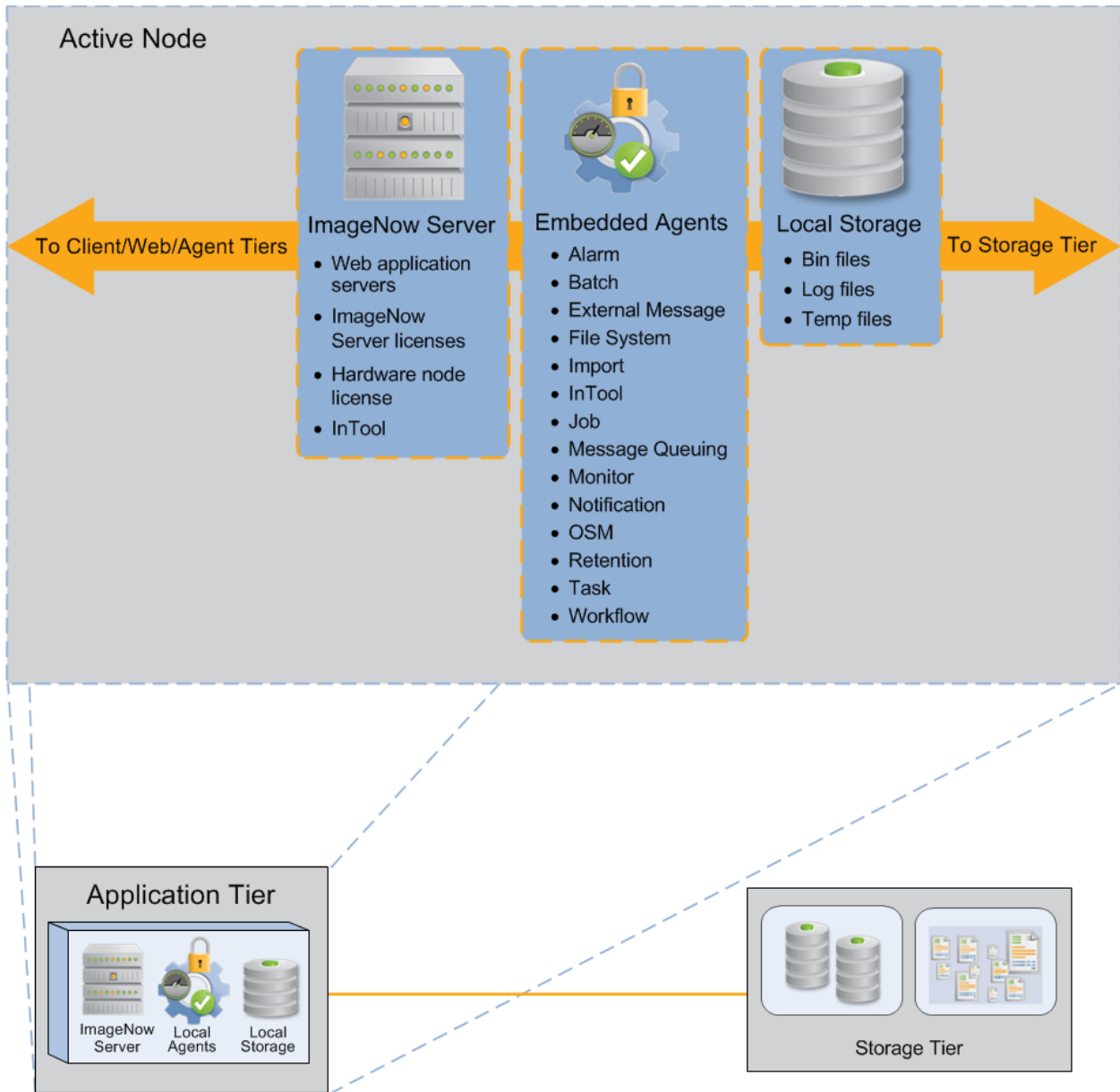
Like WebNow Server, you can combine more than one Integration Server to form a cluster, and combine two or more clusters to form a farm. Integration Server is multi-threaded, which allows for concurrent execution of multiple client requests. Integration Server supports SSL for secure client-to-server and server-to-client communication and HTTP or HTTPS transport for structured data exchange.

## Application Tier

The Application Tier contains one or more nodes where each node contains Perceptive Content Server, local configuration specifications, embedded agents, and local storage such as temp, bin, and log files. You can set up an unlimited number of nodes within a cluster, and more than one cluster to form a server farm. The embedded agents in the Perceptive Content Server process jobs and provide messaging to Perceptive Content Clients. The Perceptive Content Server directs the internal agents to perform as needed without any manual steps. However, you can manually configure some agents.

Two or more nodes always run at the same time, which allows for better utilization of failover hardware. Communications with the Client Tier and the Web Tier routes through load-balancing devices. Outgoing communications from the application tier occurs through TCP/IP and AES.
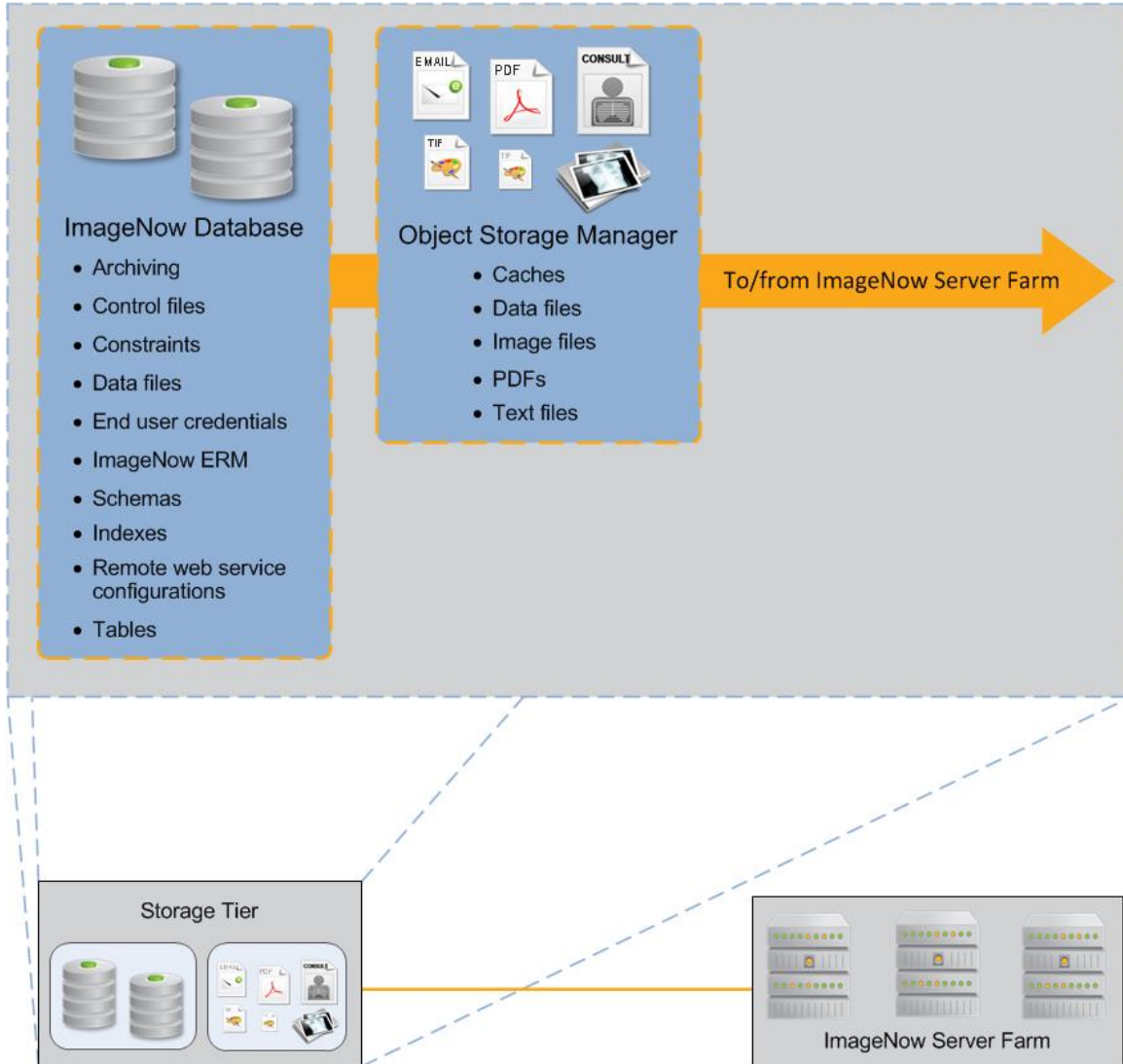
## Storage Tier

The Storage Tier contains the Perceptive Content Database and OSM, utilizes a clustered file system for the OSM, and can use block-level shared storage. The Perceptive Content Database performs all database transactions, stores metadata information related to the content stored in the OSM, as well as system information, data files, schemas, indexes, and tables.

The OSM is a tree-structure file system that consists of a main directory comprised of sets or branches. Perceptive Content stores all document objects in the OSM in their original format, for example TIFF, PDF, or Microsoft Word. You can configure the OSM to store objects across any number of file systems on a variety of platforms and architectural designs. Using INTool, a command-line tool, you can manage your OSM storage. Communication between the Storage tier and the Agent tier occurs via ODBC and SSL.

**Note** Perceptive Software recommends configuring the OSM layer in an active-passive configuration using a cluster server with active and passive OSM file servers for failover of SAN ownership to a new owner in the event of a system error.

# High availability and the Perceptive Content Server

The following section provides a high-level overview of Perceptive Content Server functionality in a high-availability configuration.

## Named server instances

During the installation of an instance of Perceptive Content Server, you must provide a unique name for the instance that is used to identify that specific instance within the system. This instance name distinguishes log and temporary file output in environments with multiple instances. You can also use this name in configuration files so that multiple instances can share the same configuration files. The instance name is often required when performing certain server management activities.

**Note**  Instance names must be unique across all instances in the server environment.

## Using references for nodes

In a highly available environment, nodes actively process data just as they do in a traditional server environment. Previously, a variable contained a path to all of the Perceptive Content Server resources and was required for the system to start and run properly. With Perceptive Content 7.0, this variable holds a path to the portion of Perceptive Content Server resources that are shared among different machines. The resources are located on an NTFS share or a mounted NFS file system.

While there can be more than one-instance-node per hardware machine, Perceptive Software recommends that there be only one instance node per machine.

## Passing references to the server

Traffic routing for a highly available Perceptive Content Server and agent connections are available for activities such as passing multiple server-IP addresses to the server. All Perceptive Content Server traffic routing occurs via TCP/IP, and WebNow traffic routing occurs using HTTP or HTTPS. Traffic routing passes through a load-balancing device at each tier level in the system.

## Using references as arguments in INTool

You use INTool commands for a highly available server just as you do for a stand-alone Perceptive Content Server configuration. You can use these commands to perform actions like creating and configuring multiple instances of Perceptive Content Server and creating and configuring alternate installation locations of Perceptive Content Server. When you use these commands for a highly available server, take care when you supply arguments that involve providing the instance name.

For more information about INTool commands, refer to topics in Administrator Help.

## Using environment settings

Some Perceptive Content Server INI file settings affect how the system operates in a high-availability environment. These settings include IMAGENOWDIR and other environment variables that appear in the environment.ini file. IMAGENOWDIR stores the path to the portion of INServer resources, including etc, OSM, form, and job files, that must be shared among different machines and must be moved to a UNC share. For more information about INI settings, refer to the *Perceptive Content Server and Client Installation and Setup Guide*.

# High availability for databases

At the database level, you must understand the following considerations to manage the increase in connections and the load on the database. Refer to your database vendor documentation for additional options for high-availability solutions, such as database mirroring.

- As the number of processes and connections rises, the amount of memory needed by the database and operating system to handle the load also increases. To ensure adequate performance, you might need to increase OS memory and the number of CPUs to support the higher demand for resources.

- For systems using an Oracle database, you may need to increase the System Global Area (SGA) memory to ensure adequate performance. For some installations of Perceptive Content, the Oracle parameter called "Processes," which is used to establish the default number of processes that are started, is set too low at 150 processes. If you add agents or increase the num.workers setting in any of the agent INI files, you must reevaluate this parameter. You can set the parameter to between 300 and 500 for performance tuning. Refer to Oracle documentation for more information about the "Process" parameter, and the "Tune Oracle for Perceptive Content Server" section in the *Perceptive Content Installation and Setup Guide* for Oracle for more information about configuring processes.

  You can also set an Oracle database to high-availability using an Oracle Real Application Clusters (RAC) database. Refer to the Oracle website and documentation for more information.

- For systems using a Microsoft SQL Server database, there are worker thread settings in the inserver.ini file that have default values. These values are optimized for the Perceptive Content Server and Client, and for most server-client configurations, do not need to be adjusted for performance. For more information about worker thread settings, refer to the inserver.ini settings table in the *Perceptive Content Server and Client Installation and Setup Guide* for Microsoft SQL Server.

# High availability for file storage systems

There are considerations you need to take into account for storage in a highly available environment, including the performance level of the device you use for traditional file system storage and third-party products you use for content-addressed storage (CAS).

## File system storage

For file system storage of OSM content, Perceptive Software recommends that you use a high performance storage-area network (SAN) device. In prior versions of this software, a Perceptive Content Server machine typically owns and mounts the SAN device. However, because of the active-active feature in the Perceptive Content 7.0 release, drives can connect using a UNC path.

## Content-addressed storage

CAS is a storage solution that retrieves data based on the content of the data and not where the data is stored. An EMC Centera storage device is one example of a CAS device. These devices typically have their own endpoint and disaster recovery solutions. For more information about CAS high-availability solutions, refer to the third-party vendor's documentation.

# High availability for WebNow

You can configure WebNow with a Network Load Balancing (NLB) clustered application server to ensure the integrity of open files for each user session and to ensure that the system distributes traffic across multiple nodes for system redundancy.

When load balancing connections for WebNow, you must set attributes in the WebNow.settings file (for example, heartbeat), to ensure that connections go to the same servlet for session affinity (also known as a "sticky" session). For more information, refer to the "Perceptive Content Server interaction" entries in the WebNow.settings properties table in the *WebNow Installation and Setup Guide*. WebNow NLB clustering includes session affinity or "sticky sessions" to ensure the integrity of open files for each user session.

# High availability for Perceptive Content Message Agent

You can set up Message Agent with multiple instances that point to a single Perceptive Content Server or to a different server. Message Agent nodes have their instance names configured and named when configuring the instances during installation or by passing in a command-line parameter.

You can configure multiple instances of Message Agent on different servers and then segment them, based on the number of users, into different groups that each has a server. Additionally, you can integrate a load balancer into the system to control traffic limitations and for additional redundancy.

# High availability for Perceptive Integration Server

To establish high availability for Perceptive Integration Server, you set up multiple instances of Integration Server to run on different servers and then segment the user base into different groups. You can then assign each group with its own server. Additionally, you can add a load balancer in front of the server for increased redundancy.

**Note**  You do not need to specify a name when creating a new instance for Perceptive Integration Server. For more information, refer to the *Perceptive Integration Server Installation Guide*.

# Index