# Continuous Operations and High Availability

## Overview

The Cisco Nexus 7000 platform is modular in its design, with an emphasis on redundant critical components throughout all of its subsystems. A highly modular and compartmentalized approach to systems design is applied to all facets of the platform, spanning across the physical, environmental, power, and system software aspects of its architecture. Additionally, a distinct functional separation between control plane and forwarding data plane is emphasized in its design in order to allow continuous operation and zero service disruption during planned or unplanned control-plane events or failures. This document is intended to provide details of the high availability features and options across all of these areas in order to provide the reader a better understanding of the hardware and software fault tolerance capabilities of the Cisco Nexus 7000 platform architecture.
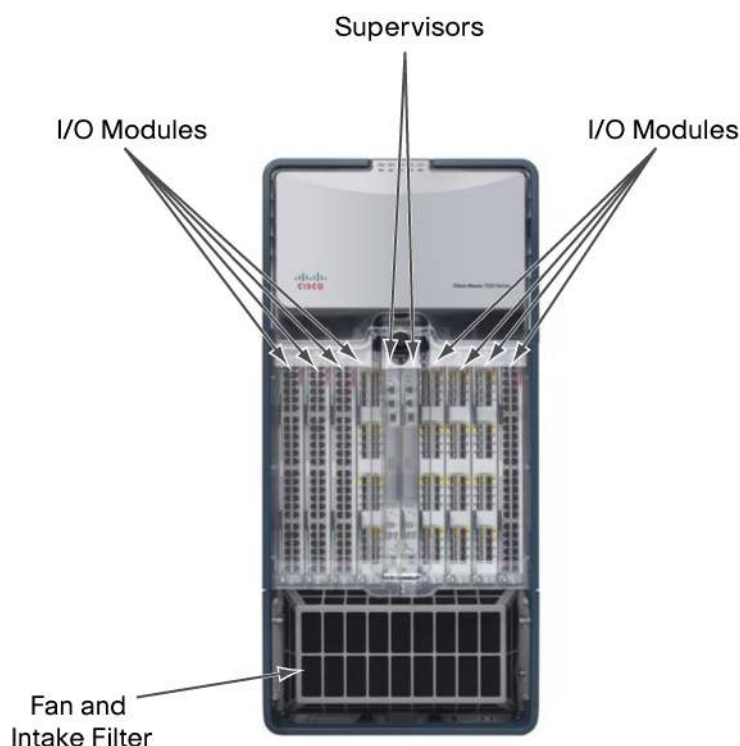
## Physical Subsystem Redundancy Options and Features

### Supervisor Module Redundancy

The Cisco Cisco Nexus 7000 platform supports dual supervisor modules to provide 1+1 redundancy for the control and management plane. A dual supervisor configuration operates in an active/standby capacity in which only one of the supervisor modules is active at any given time, while the other acts as a standby backup. State and configuration remain constantly synchronized between the two supervisor modules to provide seamless and stateful[1] switchover in the event of a supervisor module failure. NX-OS's Generic On-Line Diagnostics (GOLD) subsystem and additional monitoring processes on the supervisor facilitate the triggering of a stateful failover to the redundant supervisor upon the detection of unrecoverable critical failures, service restartability errors, kernel errors, or hardware failures.

---

[1] Stateful switchover is handled on a per service basis. Not all services have been designed to switch over in a stateful manner. See subsequent sections of this document for further details.

**Figure 1.** Cisco Nexus 7000 Series 10-Slot Chassis Front View
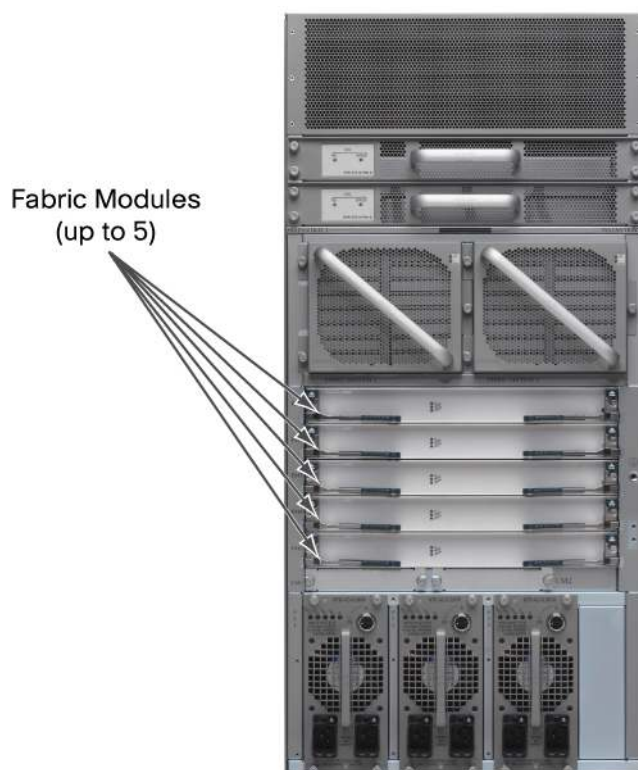


In the event of a supervisor-level unrecoverable failure, the currently active but errored supervisor will trigger a switchover. In executing the switchover, the standby supervisor becomes the new active supervisor, using the synchronized state and configuration while the errored supervisor is reloaded. If the failed supervisor is able to reload and pass self-diagnostics it will initialize, become the new standby supervisor, and synchronize its operating state with the newly active unit.

Additional details on supervisor switchover are contained in the "Supervisor Switchover" section of this document.

**Switch Fabric Redundancy**

In addition to control-plane and management availability through redundant supervisor modules, the Cisco Nexus 7000 platform also provides switching fabric availability through redundant switch fabric module implementation (Figure 2). A single Cisco Nexus 7000 chassis can be configured with one or more fabric modules, up to a maximum of five for capacity as well as redundancy. Each I/O module installed in the system will automatically connect to and use all functional installed switch fabric modules. A failure of a switch fabric module will trigger an automatic reallocation and balancing of traffic across the remaining active switch fabric modules. Replacement of the failed fabric module reverses this process. Once the replacement fabric module is inserted and online, traffic is again redistributed across all installed fabric modules, thereby restoring the redundancy level.

**Figure 2.** Cisco Nexus 7000 Series10-Slot Chassis Rear View (Fabric Modules)

Fabric Modules
(up to 5)

**Cooling Subsystem**

The Cisco Nexus 7000 10-slot chassis contains two redundant system fan trays for I/O module cooling and two additional redundant fan trays for switch fabric module cooling (Figure 3). Both types of fan tray are hot swappable to the system.

The fan speeds are variable and are dynamically adjusted to one of 16 levels[2] in order to optimize system cooling in both sections while minimizing overall system noise and power draw.
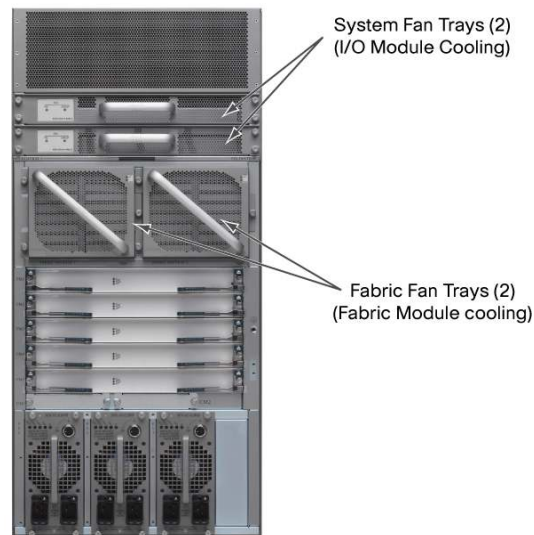
The system and fabric fan speeds are automatically adjusted dependent upon a number of factors including:

- Highest power-drawing module installed in the system or in the fabric sections
- Temperature value readings on the inlet, supervisor, and fabric module ambient temperature sensors
- Presence of the air inlet filter

Additionally, a detected failure of a fan within a given fan tray will trigger an adjustment (typically an increase) in the speed and flow rate of the remaining fans in order to compensate for the failure. Due to the thermal dynamics impact created by a missing fan tray, a detected removal of an entire fan tray, without replacement, will initiate a series of warnings followed by a system shutdown after a three-minute warning period.

---

[2] The 16 levels of variable fan output are currently a software limitation. Future versions of software may increase the granularity of output levels to up to 256 levels.

**Figure 3.**    Cisco Nexus 7000 Series 10-Slot Chassis Rear View (Fan Trays)



**Power Subsystem Availability Features**

The Cisco Nexus 7000 platform is powered by three internally redundant power supplies (Figure 4). Each of those individual power supply units is composed of two internalized isolated power units, giving it effectively two power paths per modular power supply and six paths in total, per chassis, when fully populated. The power supplies use a proportional load-sharing method for power distribution to power system components, allowing the efficient use of dissimilar capacity power supplies in the same chassis. Therefore all installed power supplies are active and share the overall load. Additionally, the power subsystem allows the three power supplies to be configured in any one of four redundancy modes.
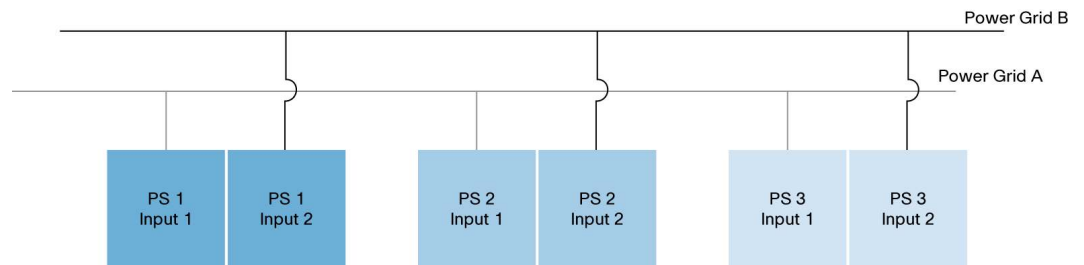
**Figure 4.**    Cisco Nexus 7000 Series10-Slot Chassis Rear View (Power Supplies)
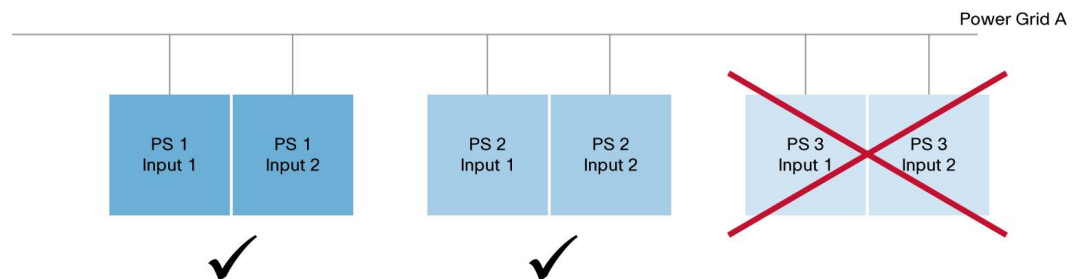
The available power supply redundancy modes are:

- Nonredundant Combined Power: In this mode, there is no power supply redundancy. All available power across installed power supplies is combined to provide the sum of all available power to the usable power budget. In this configuration, a failure of a power source or power supply will affect the overall available power to the system. If the current power draw exceeds the postfailure power budget, the system will enter a failure state.

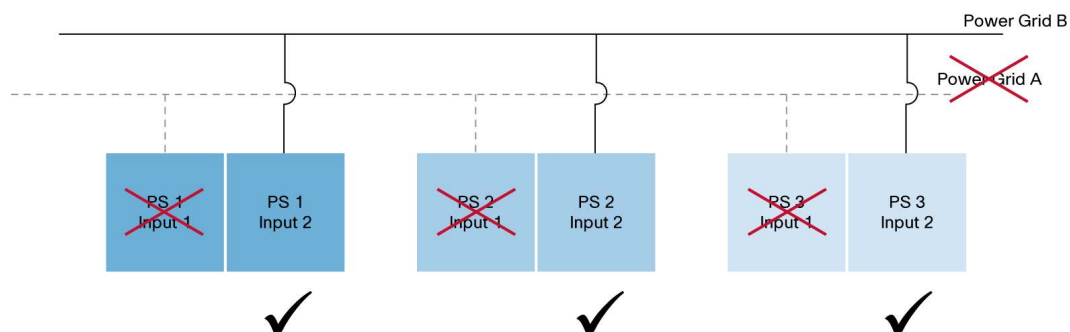**Figure 5.**   Combined power model (Non-redundant)



- N+1 Power Supply Unit Redundancy (default): In this mode, protection against a single power supply unit failure is provided. In the event of a single power supply failure, loads are redistributed using available capacity across remaining functional power supply units. N+1 redundancy is available with either 2 or 3 power supplies installed in the system.  The total available power to the chassis is the sum of all installed power supplies minus that of the largest (for redundancy).  This is the default power redundancy mode.
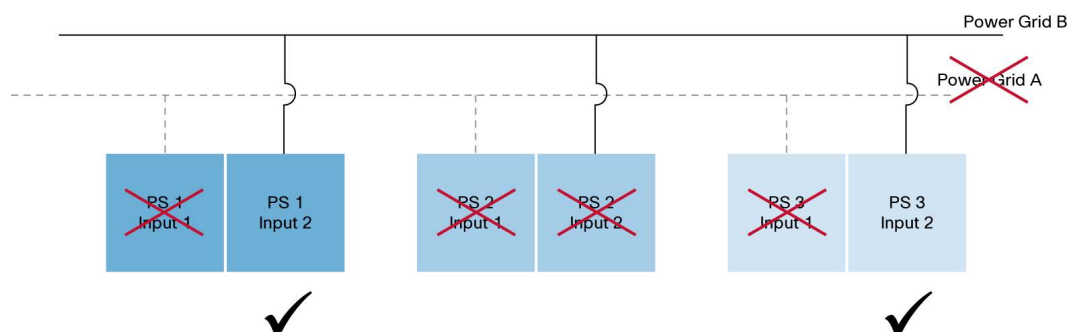
**Figure 6.**   N+1 redundancy model



- Input Source Redundancy: This mode provides protection against input source failures. In order to implement input grid redundancy, each power supply unit must be connected to two independent power sources (grids). A single power grid failure will allow all installed power supplies to continue functioning through the remaining internalized half still connected to the remaining available power grid.

**Figure 7.**   Input grid redundancy model

- Full Redundancy (Power Supply + Input Grid): In full redundancy mode, the power supply configuration provides protection against either a single power source grid failure or a single power supply unit failure. It is a logical combination of the N+1 and input grid redundancy modes.

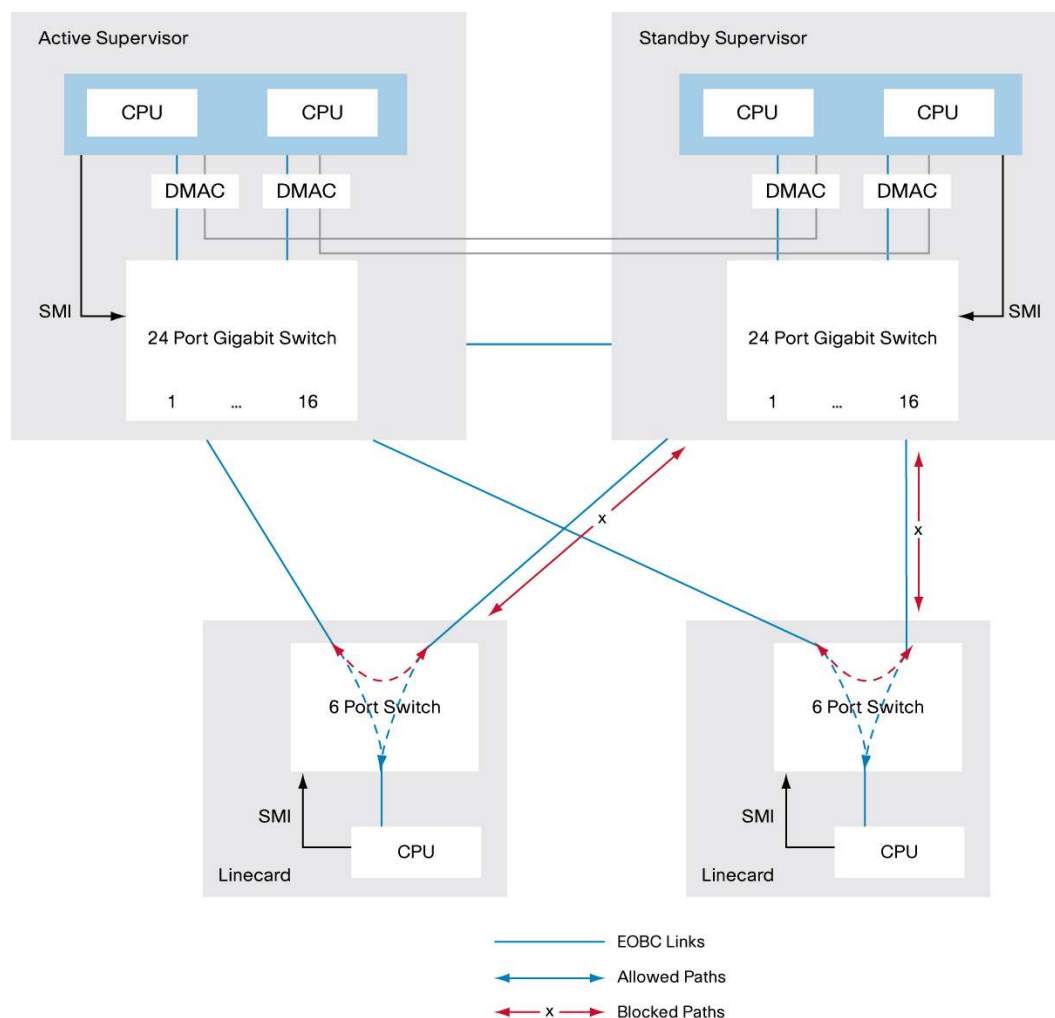**Figure 8.** Full redundancy model



Each of the power supply redundancy modes imposes different power budgeting and allocation models, which in turn deliver varying usable power yields and capacities. For more information or details regarding power budgeting, usable capacity, and planning requirements please refer to the Cisco Nexus 7000 environmental data sheet or white paper.

**Internal Ethernet Out-of-Band Channel**

Cisco Nexus 7000 uses a switched Ethernet out-of-band channel (EOBC) for management and control traffic between the supervisors and I/O modules and between the supervisors and fabric modules. On the supervisor modules, Ethernet connectivity is provided using an on-board 24-port Ethernet switch on a chip, with a one 1 Gbps Ethernet link from each supervisor to each I/O modules, each supervisor to each switch fabric module (up to five), and between the two supervisors (Figure 9). Two additional redundant 1 Gbps Ethernet links are used on each supervisor to connect to the local CPU within the supervisor. This design provides a highly redundant switched-Ethernet-based fabric for control and management traffic between the supervisors and all other processors and modules within the system.

**Figure 9.** EOBC Connectivity in Cisco Nexus 7000

The modules utilize a six-port Ethernet switch on a chip, of which only three ports are used to provide three 1 Gbps lines of connectivity. One line is used for connectivity to the module local CPU, while two 1 Gbps lines are used to connect to each of the supervisors. These Ethernet-based channels are used strictly for control-plane signaling and do not carry data-plane traffic.

## NX-OS System Software Redundancy Options and Features

The Cisco Nexus 7000 platform runs the NX-OS operating system. NX-OS uses a highly modularized architecture that compartmentalizes components for redundancy, fault isolation, and resource efficiency. Functional feature components operate as independent processes referred to as services. NX-OS services implement availability features by design into each service, as needed. Most system services are capable of performing stateful restarts, thereby allowing a given service experiencing a failure to be restarted and to resume operations transparently to other services within the platform, and fully transparently to neighboring devices within the network. Back-end management and orchestration of processes, services, and applications within a platform are handled by a set of high-level system-control services: the system manager, persistent storage services; message and transaction services; and the Installer.

These high availability infrastructure components provide the services, APIs, monitoring, and control support that facilitate the service restart and supervisor switchover capabilities of the platform. By employing multiple levels of service and component monitoring, combined with various layers of structured failure scenario handling, the NX-OS software architecture provides a

very comprehensive approach to helping ensure system availability. Table 1 illustrates the various potential failure scenarios and handling capabilities built into NX-OS.

**Table 1.**    Various failure scenarios handled by System Manager

| Failure Scenario | Action |
|---|---|
| Service/process exception | Service restart |
| Service/process crash | Service restart |
| Unresponsive service/process | Service restart |
| Repeated service failure | Supervisor reset (single)/switchover (dual) |
| Unresponsive system manager | Supervisor reset (single)/switchover (dual) |
| Supervisor hardware failure | Supervisor reset (single)/switchover (dual) |
| Kernel Crash | Supervisor reset(single)/switchover (dual) |
| Unresponsive System | Supervisor reset(single)/switchover (dual) |

The following sections of the document examine the details of the components that deliver these continuous operations and failure recovery capabilities.
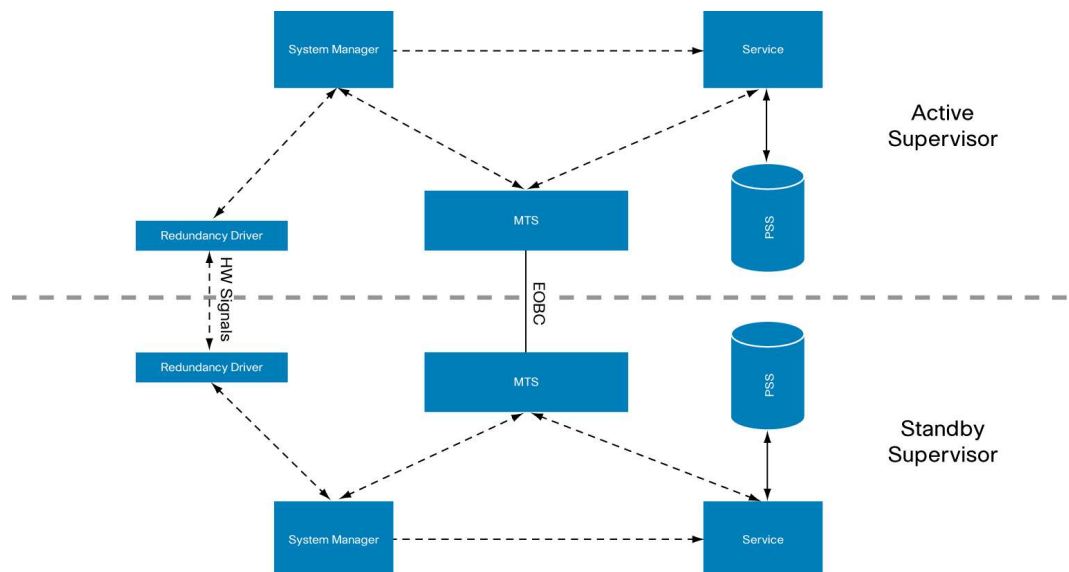
**System Redundancy Components**

NX-OS leverages three key core infrastructure services to provide overall high availability functionality:

- System manager

- Persistent storage service

- Message and transaction service

In a redundant configuration, such as when dual supervisor modules are in operation, mirrored services run on each supervisor module, with configuration and operating state synchronized between them. One of those supervisors will operate as the active supervisor while the other operates in a standby facility until activated in a switchover.

**Figure 10.**    NX-OS High-Availability core infrastructure components
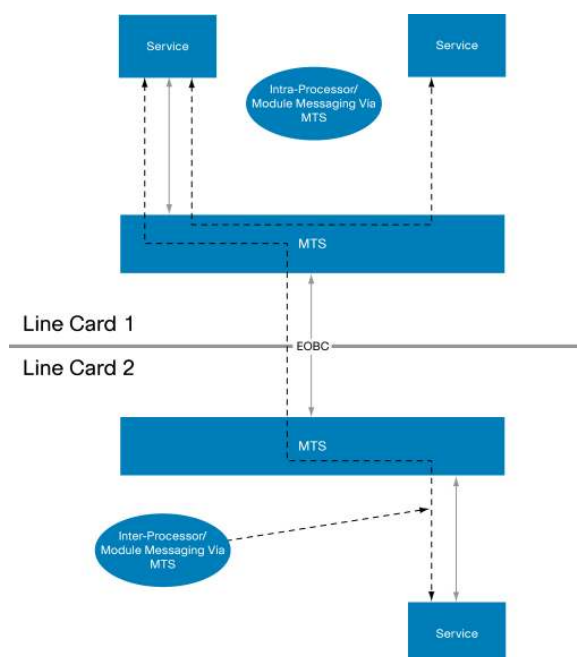
The system manager orchestrates overall system function, service management, and system health monitoring. It is also responsible for maintaining overall high availability states, enforcing high availability policies, and managing system configuration. It is the system manager that is responsible for launching, stopping, monitoring, and restarting services. System manager is is also leveraged for initiating and managing the syncing of service states and intersupervisor states for stateful switchover. It is also the system manager that will initiate a supervisor switchover if it determines that the current supervisor has undergone an unrecoverable failure or if key core services are undergoing errors and cannot be restarted reliably. In addition, being the overall control and monitoring process, the system manager is responsible for "punching" or triggering the keep-alive indicator for the hardware-based watchdog timer on the supervisor. The lack of this periodic heartbeat from the system manager within the keep-alive timeout period of the watchdog timer will indicate a nonresponsive system manager, which will trigger a hardware-based supervisor reset (single supervisor) or switchover (dual supervisors). The system manager's health is also monitored by a kernel level module receiving periodic heartbeats sent by the system manager process. This allows the system to take corrective action in response to an unresponsive system manager that has exceeded the heartbeat timeout period.

The persistent storage service (PSS) is the base infrastructure component responsible for storing and managing the operational run-time information and configuration of the other platform services. It is leveraged by other system services to recover state in the event of a service restart. PSS provides a managed, structured API to read and write data to and from the storage system, essentially functioning as a database of state and run-time information. Services wishing to leverage the PSS infrastructure are able to checkpoint their state information periodically, as needed. This allows services to subsequently recover to the last known operating state preceding a failure, thereby allowing for a stateful restart. This state recovery capability is available to NX-OS services in both single- and dual-supervisor configurations, and helps enable a service to transparently return to operation without service disruption to data-plane traffic, neighboring devices, or other internal services.

For example, even in a single supervisor configuration, the PSS enables the stateful restart of a service such as spanning-tree without impacting the overall spanning-tree topology or stability.

The NX-OS message and transaction service (MTS) is a high performance interprocess communications (IPC) message broker that specializes in high availability semantics. MTS handles message routing and queuing between services on and across modules (including across supervisors). MTS facilitates the exchange of messages such as event notification, synchronization, and message persistency between system services and system components. MTS also facilitates and manages message queuing between services and, as a result, can maintain persistent messages and logged messages in queues for access even after a service restart.

**Figure 11.**   MTS inter-service communications

Each of these system infrastructure component services plays a role in providing overall system availability. Individual functional system services (for example, authentication, authorization, and accounting [AAA], Spanning Tree Protocol, Network Time Protocol [NTP], command parser, and so on) are responsible for using these system-level redundancy infrastructure components to checkpoint their own states, recover the states upon restart, and synchronize the states to their standby counterparts using the PSS and MTS facilities.

**Service Modularity and Restart Capabilities**
The services within NX-OS are designed as nonkernel space (user space) processes that perform a function or set of functions for a subsystem or feature set.

Each service (essentially a feature) and each service instance is run as a separate independent protected process. This approach provides a highly fault tolerant software infrastructure and fault isolation between services. In short, a failure in a service instance (such as 802.1q or MST) will not affect any other services running at that time (such as LACP). Additionally, each instance of a service can run as an independent process. This implies that two instances of a routing protocol (for example, two instance of OSPF) run as separate processes, thereby providing fault isolation even between those two instances of the same service.

This resilient highly modular architecture allows the system to provide the highest level of protection and fault tolerance for all services running on the platform. This approach also facilitates rapid fault isolation, resolution, and modular software upgrades to address specific issues, while minimizing the impact to other critical services and the overall system.
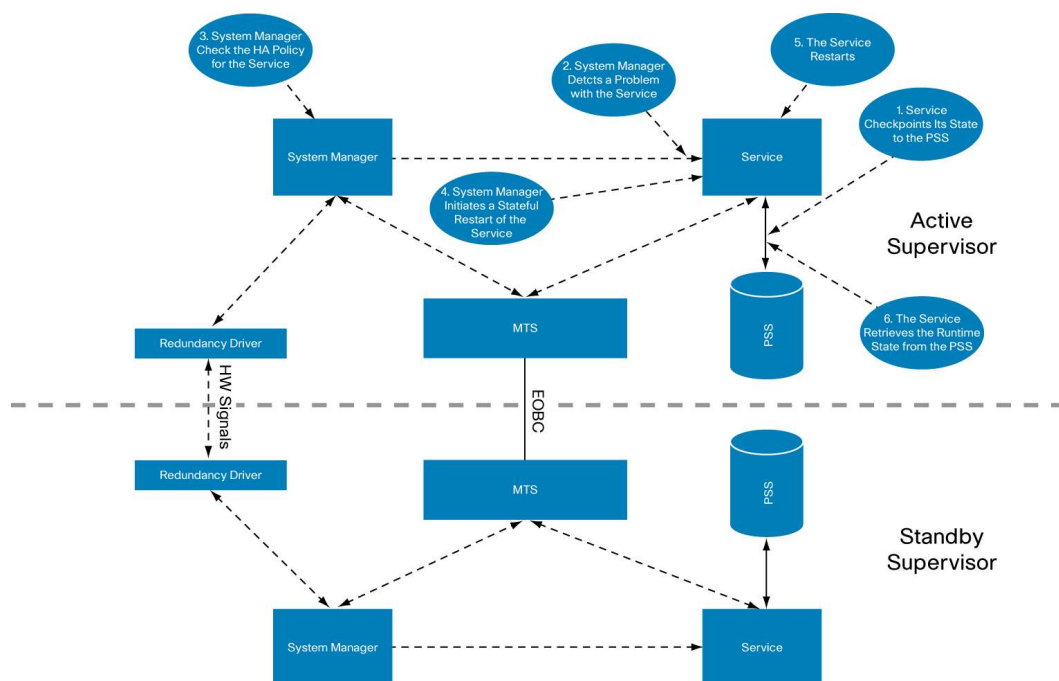
The services in NX-OS are also capable of undergoing rapid restart. A service restart may be initiated automatically by the system manager in the event of critical fault detection. When a restart is initiated, a service process is sent a signal to stop and is subsequently shut down and then rapidly restarted, typically in the order of milliseconds. This allows an error condition to be cleared and a service to be reset if necessary.

A service can undergo different types of restarts, stateful or stateless. A service may be restarted by the system manager depending on current errors, failure circumstances, and configured high

availability policy for the service. If the service is issued a stateful restart, the new service process will retrieve the previous run-time state data from PSS and resume operations from the last checkpointed service state. Most of the services in NX-OS are designed to support stateful restart capabilities by leveraging the high availability infrastructure services of PSS and MTS. If the service is issued a stateless restart, it will initialize and run as if it had just been started with no prior state.

Not all services are designed for stateful restart. More specifically, currently the group of Layer 3 routing protocols (for example, Intermediate System-to-Intermediate System [IS-IS], OSPF, BGP, Routing Information Protocol [RIP], and so on) for IPv4, IPv6, and IP multicast are not designed to leverage the state persistence of PSS within NX-OS. Instead, these protocols are currently designed to rebuild their operational state using information obtained from neighbors. Additional details on the high availability functionality of Layer 3 protocols are addressed in the "Layer 3 Protocol High Availability Features" section of this document.

**Figure 12.**   Service restart process



### Supervisor Switchover

The CISCO NEXUS 7000 platform provides 1+1 redundant supervisor modules that can perform a supervisor switchover (SSO) in critical failure situations. The switchover of active to standby supervisor is nondisruptive to the forwarding plane and can therefore provide nonstop forwarding of data during a supervisor switchover.

The active and standby supervisor modules in the Cisco Nexus 7000 platform remain synchronized in state and configuration using a two-stage process. The first time a supervisor enters a fully operational state, it will determine whether it is the active or the standby supervisor. If it is the standby unit, it will perform the first stage of the synchronization by initiating a global sync (GSYNC) with its counterpart. To accomplish this, the standby supervisor will issue a GSYNC request to the active supervisor. The active supervisor unit, in return, will begin a complete exchange of current operating configuration and state to the standby unit. Once this first-stage complete information exchange has taken place, the supervisors enter the second stage of SSO synchronization, in which event driven and state change triggered sync updates are sent through messages through the MTS facility. A service that is updating or checkpointing its state will issue a message to MTS as well as writing to PSS. MTS will route the subsequent sync message to the peer MTS process on the standby supervisor through MTS messages over the Ethernet out-of-band channel. The MTS process on the standby supervisor will notify the corresponding process on the standby supervisor in order for it to update and checkpoint its own state to match that of the peer process on the active supervisor.

Additionally, the system manager of the active supervisor will monitor the state and operational readiness of the standby supervisor. If the standby supervisor has entered an unexpected state, or if MTS synchronization to the standby fails, the standby supervisor is reset in an effort to facilitate effective and consistent state between the active and standby supervisors.

In the event of the detection of an unrecoverable error, the system manager, depending on the configured high availability policy, may initiate a supervisor switchover from the current active supervisor over to the current standby unit. A switchover may be triggered based on one or more criteria that contribute to the determinations of an unrecoverable failure:

- Repeated service restarts within the "restart timer" stability window
- Hardware failure
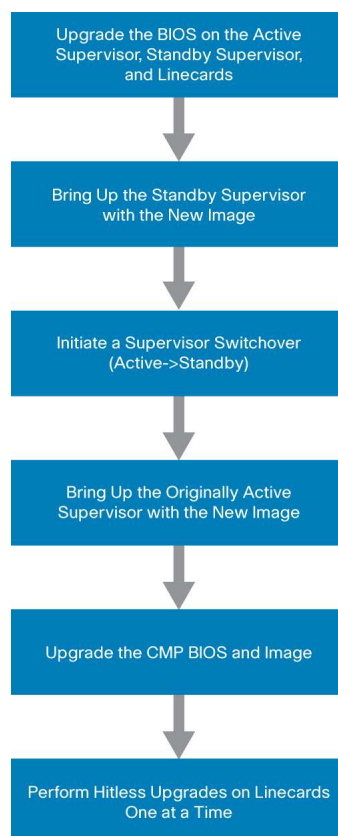- Kernel panic

**In-Service Software Updates**

NX-OS provides the ability to perform in-service software updates (ISSUs). The overall modular software architecture of NX-OS supports plug-in-based services and features. This framework makes it possible to perform complete image upgrades, nondisruptively without impacting the data forwarding plane. This capability allows for nonstop forwarding during a software upgrade, including upgrades between full image versions (for example, from 4.0 to 4.1).

An in-service software upgrade is initiated manually either through the command-line interface (CLI) by an administrator, or (in future releases) via the management interface of the Datacenter Network Manager software platform. When initiated, an in-service software upgrade updates (as needed) the following components on the system:

- Supervisor BIOS, Kickstart Image, System Image
- I/O module BIOS and Image
- Connectivity Management Processor (CMP) BIOS and Image

Once initiated, the ISSU Installer service begins the ISSU cycle. The upgrade process is composed of several phased stages designed to minimize overall system impact with no impact to data traffic forwarding. Figure 13 provides a high-level illustration of the in-service software upgrade process.

**Figure 13.**   In-Service Software Upgrade Process



It is important to note that an ISSU-based upgrade is a systemwide upgrade and not a virtual device context (VDC)–based upgrade. Therefore, the initiation of an ISSU-based upgrade applies the same image and versions across the entire system, to all configured VDCs. Virtual device contexts are primarily a control-plane and user-interface virtualization and cannot currently run independent image versions per virtualized resource.

Additionally, in order to minimize overall network impact, an in-service software upgrade should only be initiated when the network topology is stable, as control-plane components necessary to participate in topology state changes may be in the process of upgrading during a transition.

**Layer 3 Protocol High Availability Features**

The Layer 3 protocols for IPv4, IPv6, and IP multicast do not currently implement stateful restart capabilities. This capability may be added in the future for these protocol services. Currently, the Layer 3 protocols utilize two major methods for operational recovery during a restart:

- Graceful restart extensions
- Protocol-based periodic refresh

OSPFv2, OSPFv3, IS-IS, Enhanced Interior Gateway Routing Protocol (EIGRP), and BGP utilize graceful restart extensions to the base protocols to provide nonstop forwarding and least obtrusive routing recovery for those environments. The NX-OS routing protocol extensions for graceful restart follow the standards outlined in RFCs 3623, 4724, and 3847 for OSPFv2/v3, BGPv4, and IS-IS respectively. It is important to note that while the NX-OS standards-based implementations of graceful restart extensions and non-stop forwarding are compatible with Cisco IOS versions which also leverage the same IETF standards, they are not compatible with older IOS versions that use the pre-standard implementations of Cisco Non-Stop Forwarding (NSF).

In addition to the standards-based implementations for the aforementioned protocols, graceful restart extensions have been developed for use with EIGRP. The extensions to EIGRP used in NX-OS are compatible with those used for EIGRP nonstop forwarding on other Cisco platforms.

Graceful restart allows the Cisco Nexus 7000 platform to maximize the benefits of the distinct separation between control plane and data forwarding plane in its architecture. If a restart of a graceful-operations-capable protocol is required, that particular routing protocol will utilize the relevant mechanisms to signal to its neighbors that a planned restart is being executed. This notification to its neighbors will allow the neighboring devices to continue forwarding traffic to the restarting entity as if operations were normal and allows the restarting service to remove itself from the network control plane non-disruptively. The unit restarting its routing protocol will continue to forward traffic based on the last established routing and forwarding information bases (RIB/FIB), independent of control-plane operations, thereby allowing uninterrupted, continuous data delivery during the restart. Once the restarted routing service has reestablished a stable state, it can notify its neighbors and rebuild its adjacencies, thereby nondisruptively reinserting itself into the network from a control-plane perspective.

The RIPv2, Protocol Independent Multicast (PIM), PIM6, Internet Group Management Protocol (IGMP), MSDP, and MLD protocols do not implement graceful restart extensions. Instead, as refresh-based protocols, they utilize the periodic refresh from neighbors that is inherent to the protocol to reestablish their state.

Since the routing protocols themselves are not stateful during recovery from a restart, neither are the routing information bases for these protocols. The unicast and multicast RIBs must be reestablished and rebuilt based on the recovered state postrestart of the relevant Layer 3 protocol. RIB and FIB are maintained during the restart in order to continue forwarding operations, but are updated once the protocol is stable and has reestablished its recovered state. This allows for any necessary updates to RIB and FIB, as needed, to accommodate any network changes after restart.

**Virtual Device Contexts**

NX-OS implements a logical virtualization at the device level allowing multiple instances of the device to operate on the same physical switch simultaneously. These logical operating environments are known as virtual device contexts, or VDCs. Virtual device contexts provide logically separate device environments that can be independently configured and managed. This degree of isolation between virtual devices provides a very valuable level of fault isolation in addition to the security and administrative benefits. Human error or failure conditions due to configuration are naturally isolated within a given virtual device. While virtual device contexts are not primarily a high availability feature, their effective creation of operationally independent fault

domains, contributing to availability and preventing service disruption associated with device configuration, should always be considered.

The specific functional and operational details of virtual device contexts are outside the scope of this document. For more information on VDCs, please refer to the white paper on Cisco Nexus 7000 virtualization features.

## Summary

With this deeper examination of the Cisco Nexus 7000 platform, it becomes apparent that the entire platform is truly designed from the bottom up with a goal towards highly-available uninterrupted service under almost any condition. By strategically employing fault isolation, fault detection, and fault recovery on many levels, in all areas, the Cisco Nexus 7000 platform is capable of delivering critical enterprise-class service for current and future business needs while preserving capital investment as well as the data path.

**Glossary of Terms Used**

NX-OS: Datacenter Operating System

MTS: Message and Transaction Service

PSS: Persistent Storage Service

EOBC: Ethernet Out-of-Band Channel

IETF: Internet Engineering Task Force

RFC: Request For Comment

CMP: Connectivity Management Processor

ISSU: In-Service Software Update

RIB: Routing Information Base

FIB: Forwarding Information Base

VDC: Virtual Device Context

NSF: Non-Stop Forwarding

HA: High-Availability

API: Application Programming Interface

MST: Multiple Spanning-Tree

LACP: Link Aggregation Control Protocol

Printed in USA

C11-446927-01   02/08