



---

SPARC SERVERS

An Oracle White Paper  
March 2013

# Maximizing Application Reliability and Availability with the SPARC M5-32 Server

Introduction .....	2
Reliability, Availability, and Serviceability (RAS) .....	3
Trade-offs .....	5
RAS from the Inside Out.....	5
SPARC M5 Processor .....	5
SPARC M5-32 System Memory .....	6
SPARC M5-32 Server .....	6
SPARC M5-32 Hypervisor .....	6
SPARC M5-32 I/O .....	7
Platform Management .....	8
Oracle Solaris.....	8
Oracle VM Server for SPARC.....	10
Oracle Solaris Cluster.....	10
Integrated System Monitoring.....	11
Oracle Enterprise Manager Ops Center.....	12
Oracle Auto Service Request.....	12
Conclusion .....	13

## Introduction

Oracle's SPARC M5-32 server is a high-performance system that has been designed, tested, and integrated to run a wide array of enterprise applications. It is well suited for Web, database, and application components. This versatility, along with powerful, bundled virtualization capabilities, makes it an ideal platform on which to consolidate large numbers of applications, databases, and middleware workloads or to deploy complex, multiuser development, test, and deployment environments.

The reliability, availability, and serviceability (RAS) characteristics of the SPARC M5-32 server go well beyond its highly reliable components. Redundant hardware is combined with software RAS capabilities and advanced and integrated management and monitoring. Together, these technologies enable the SPARC M5-32 server to deliver mission-critical uptime and reliability. The SPARC M5-32 server architecture enables a high degree of isolation between concurrently deployed applications, which might have varied security, reliability, and performance requirements.

The SPARC M5-32 server provides an optimal solution for all workloads, ranging from batch processing, to data warehouse applications, to highly concurrent online transaction processing (OLTP) applications. With its combination of high core and thread count, and up to 1 TB of memory per processor, the SPARC M5-32 server delivers extreme capacity and performance in a highly available, highly secure environment.



Figure 1. SPARC M5-32 server.

## Reliability, Availability, and Serviceability (RAS)

In casual conversations about system and services availability, the term *RAS* is often used to reflect availability. However, the other two components of RAS—reliability and serviceability—are equally important when considering the capabilities of a delivery platform or data center solution.

- **Reliability**

There are several ways to deliver reliability. The first method is to reduce the number of components in a system. This reduction of components drives higher mean time between failures (MTBF).

A second method for better reliability is to reduce the number of connectors for ASICs and boards. Larger servers tend to have more hot-swappable or hot-pluggable components than rack form factor servers. This means far more connectors, which lowers reliability.

The third method for better reliability is data integrity. Features such as error-correction code (ECC), data parity, and instruction-level retry deliver this. It might be assumed that large servers typically used at the back-end layer of the data center's tiered compute model are selected because they are more reliable. Large servers have a higher component count because they are designed for higher availability.

- **Availability**

Availability is provided through three key methods. The first is redundancy of components. This adds cost, but it is easily justified when the cost of an outage is greater than the cost of redundant components.

A second method for greater availability focuses on data availability. This is delivered through Extended-ECC memory (a spare dynamic random access memory [DRAM] on every dual inline memory module [DIMM]), memory lane sparing, memory pin steering, processor interconnect lane sparing, and the live migration of virtual machines between servers. Data availability is also provided by Oracle's ZFS file system and implementing mirroring across disk drives.

A third feature that delivers greater availability focuses at the data services level: clustering. A balance has to be made between redundant components in a server that still contains single points of failure and a clustered solution that focuses on data service availability by constantly monitoring the health of the servers, storage, and networking, as well as the applications.

- **Serviceability**

Components can be either hot-swappable or hot-pluggable to be counted as serviceable while the system is powered on and active. Hot-swapping of components is the least disruptive RAS feature because it requires no assistance or preparation of the host operating system and server.

Components can simply be replaced while the system continues to run, with no setup required.

Another serviceability feature is hot-plugging, which is similar to hot-swapping in that it also allows the guest operating system to continue running while the component is inserted or removed. However, a hot-pluggable component requires that the operating system and platform be prepared before the component is removed or inserted. Hot-pluggable components sometimes also require additional setup commands in order to properly use the device. Disk drives are an example of a hot-pluggable device.

Hot plugging and hot swapping are a necessity of large servers with redundant components, but they are used far less with clustered solutions. For SPARC M5-32 servers configured for redundancy, individual components such as service processors or PCIe cards can be serviced while the system is running.

High levels of RAS are achieved through a combination of several system elements:

- Processor
- Memory
- System I/O
- Service processor
- System ASICs
- System power and cooling
- Oracle Solaris operating system
- Oracle Solaris Cluster software

Each RAS element builds on top of the other to deliver the overall RAS capability within a complete system:

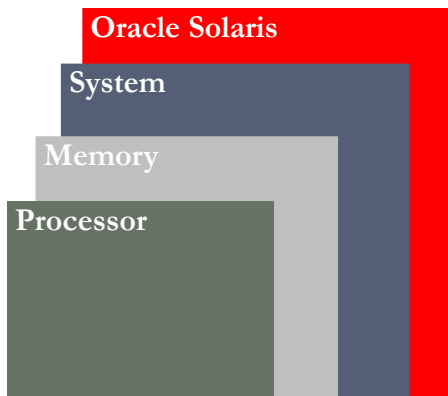


Figure 2. Elements that contribute to RAS.

This layering of system elements provides a complete RAS picture for the SPARC M5-32 server's compute nodes.

If the intent is to only deploy application services on a single server, then this stack is sufficient. For even higher level of RAS support, Oracle Solaris Cluster is designed to provide availability for the applications and to provide Oracle Real Application Clusters (Oracle RAC) for database availability.

The software layers also impact RAS. First and foremost is Oracle Solaris. The SPARC M5-32 server supports both Oracle Solaris 11 and Oracle Solaris 10. The Predictive Self Healing, Fault Management Architecture, Service Management Facility features of Oracle Solaris contribute greatly to the availability of hardware and application data services. This is described in greater detail later.

## Trade-offs

It is important to remember that some features delivering RAS are a trade-off with other features. For example, integrating multiple features into a single chip, or surface-mounting all chips and ASICs to a single motherboard, delivers higher reliability, but causes reduced availability (due to mean time to repair) and reduced serviceability. Another example is hot-swappable or hot-pluggable components, which require additional connectors for the components (among other things). While these features improve availability, they reduce reliability. Selection of such features is always a business decision, since there must be a balance between system cost and data services availability.

## RAS from the Inside Out

The heart of the RAS features in the SPARC M5-32 server is the SPARC M5 processor. The SPARC M5 processor is a member of the S3 core family of processors that include the SPARC T4 and SPARC T5 from Oracle.

### SPARC M5 Processor

The SPARC M5 processor is a 28 nm processor featuring six S3 cores, 128 threads, dedicated 128K L2 cache per core, 48-MB shared L3 cache, and a high frequency (3.6 GHz) clock. The SPARC M5 processor has a new high-performance directory-based protocol, which enables SPARC M5 systems to scale directly to eight sockets without any additional silicon needed to maintain memory coherence links.

Architected to reduce planned and unplanned downtime, the SPARC M5 processor includes advanced reliability and availability capabilities to avoid outages and reduce recovery time. Design features that boost the reliability of SPARC M5 processor include the following:

- L1\$ tag parity protection, along with retry on error
- L2\$/L3\$ single error correction/double error detection (SEC/DED), along with cache-line sparing to improve availability
- L2\$/L3\$ status and directory SEC/DEC protection
- Hypervisor correction and retry on L2\$ for SEC/DED events

## SPARC M5-32 System Memory

Because the large number of cores and memory make the SPARC M5-32 server ideal for consolidation, memory integrity is critical for system availability. The DDR3 1066-MHz memory in every SPARC M5-32 server has several features to improve both reliability and availability:

- All memory has ECC protection and correction.
- With Extended ECC, every DIMM has a spare DRAM that takes over if a different DRAM fails. This failover is transparent to applications.
- Each of the four built-in memory controllers in the SPARC M5 processor has memory lane sparing. In the event that a particular lane fails, the memory controller can continue without interruption to access the main memory.

## SPARC M5-32 Server

There are several elements of the SPARC M5-32 server that aid in overall reliability and availability. These include the interconnect between processors, fault diagnostics, power and cooling.

- The SPARC M5-32 supports Dynamic Domains, also referred to as Physical Domains (PDoms). These PDoms have full hardware and software fault isolation. A fault in one domain has no effect on a different domain.
- The interconnect on the SPARC M5-32 server has lane sparing. This means there are no lost packets of data during cache coherency or for remote memory access.
- The scalability switches used to scale beyond eight sockets have redundancy built-in. The system will recover if any switch board fails.
- The lanes between all scalability switches also utilize lane sparing, so no data is lost if a lane fails.
- There are redundant clock boards, so the SPARC M5-32 server can recover from a complete clock board failure.
- There are redundant service processors (SP) that will automatically failover if an SP fails.
- The SP boards are hot-pluggable for online servicing.
- All diagnostics are to the field replaceable unit (FRU) level on first fault detections.
- Advanced Power Management allows the administrator to set power policies to closely reflect data center policies. This includes the setting of power caps to keep power utilization under control, as well as to help reduce cooling requirements.
- All fans and power supplies are hot-swappable.
- Each system is dual power grid ready.

## SPARC M5-32 Hypervisor

The hypervisor firmware is in every SPARC M5-32 server, and it enables logical domain partitions for software fault containment. Logical domains are the primary virtualization technology of the SPARC

M5-32 server. The hypervisor provides processor support for error clearing, correction, and collection. There is a hypervisor installed in every Dynamic Domain (i.e. PDom). The SPARC M5-32 server supports up to four PDom.

## SPARC M5-32 I/O

The SPARC M5-32 server provides high-performing PCIe Gen 3 connectivity. The SPARC M5-32 server is among the first UNIX platforms to offer full PCIe Gen 3 functionality. Every PCIe slot utilizes a PCIe carrier that accepts low-profile PCIe cards, which can be up to x16 mechanically, but will run at a maximum of x8 electrically. The combination of 16 PCIe slots, eight disk drives, and four base I/O cards is called the I/O Unit (IOU). The IOU is logically divided in to two separate I/O subsystems. There are four IOU boards in the SPARC M5-32 server. Figure 3 shows the layout.

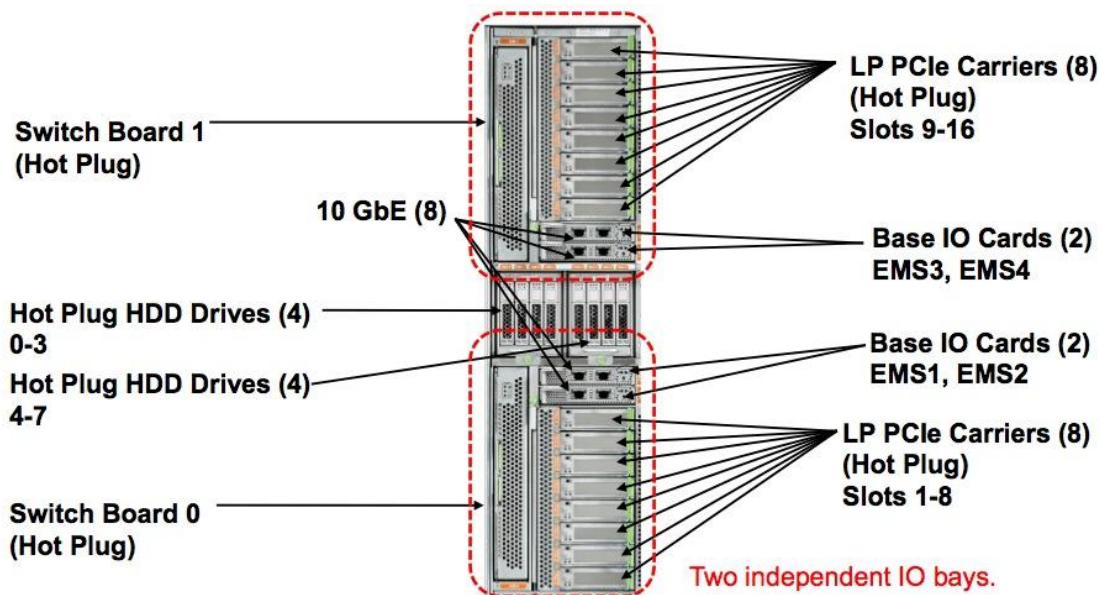


Figure 3. Layout of an IOU.

The reliability and availability features of the I/O subsystem are

- PCIe link retry is done on single-bit errors.
- All PCIe carriers can be hot-plugged while the system is running. Additional networking or storage cards can be added to a system without shutting down services.
- With redundant I/O switch boards in each IOU, the use of multipathing for networking and storage connectivity can be across the redundant switches.
- All disk drives support hot-plugging for ease of servicing and adding capacity.
- All disk drives are dual ported. This means each disk drive can have a redundant base I/O card so the drives are always accessible.



## Platform Management

Platform management is done via the service processor, or Oracle Integrated Lights Out Manager (Oracle ILOM 3.2). Oracle ILOM provides a command-line interface (CLI), a Web-based graphical user interface (GUI), and Intelligent Platform Management Interface (IPMI) functionality to aid out-of-band monitoring and administration. The management software, Oracle Enterprise Manager Ops Center, communicates with Oracle ILOM to manage and monitor the SPARC M5-32 server. All system telemetry and health diagnostics are recorded by Oracle ILOM and forwarded to Oracle Enterprise Manager Ops Center for further analysis and action. If a service event is deemed necessary, Oracle Enterprise Manager Ops Center works with Oracle [Auto Service Request](#) to notify the service that action needs to be taken.

## Oracle Solaris

The SPARC M5-32 server supports both Oracle Solaris 10 and Oracle Solaris 11. Oracle Solaris 11.1 is required in the control domain or for bare-metal operations. In logical domain guests, Oracle Solaris 10 9/10 or later is supported, with the appropriate level of patches.

Oracle Solaris includes the following features that strengthen the reliability and availability of the SPARC M5-32 server:

- **Fault Management Architecture**—The Fault Management Architecture feature of Oracle Solaris reduces complexity by automatically diagnosing faults in the system and initiating self-healing actions to help prevent service interruptions. This software helps increase availability by configuring problem components out of a system before a failure occurs—and in the event of a failure, this feature initiates automatic recovery and application restart using the Oracle Solaris Service Management Facility feature.

The Fault Management Architecture diagnosis engine produces a fault diagnosis once discernible patterns are observed from a stream of incoming errors. Following diagnosis, the Fault Management Architecture provides fault information to agents that know how to respond to specific faults. While similar technology exists with other hardware and software vendors, each of their solutions is limited to either software- or hardware-based fault detection. Oracle Solaris provides a fault detection and management environment that allows full integration of hardware fault messages to be passed on to the operating system, and the operating system services can adjust hardware resources as needed so outages are greatly reduced.

With Oracle ILOM 3.2 running on the SPARC M5-32 server's Service Processor (SP), the SP supports generation of Fault Management Architecture events to alert the administrator of platform faults. Since the SP has a simplified diagnostic engine, a running Oracle Solaris domain is not required when generating Fault Management Architecture events. The SP sends the events to the control domain, so the Oracle Solaris guest can take action on processor or memory faults. Oracle Solaris will offline cores, threads, or memory pages that generate an inordinate number of correctable errors. Taking resources offline before they generate a fatal fault helps ensure increased application service uptime. The integration of Fault Management Architecture on the SP is an availability advantage of the SPARC M5-32 server over comparable platforms from other vendors.

- **Service Management Facility**—With the Service Management Facility of Oracle Solaris, system administrators can use simple command-line utilities to easily identify, observe, and manage the services provided by the system and the system itself. The Service Management Facility describes the conditions under which failed services may be automatically restarted. These services can then be automatically restarted if an administrator accidentally terminates them, if they are aborted as the result of a software programming error, or if they are interrupted by an underlying hardware problem. Other operating systems on the market today use either a monolithic startup script or a series of smaller scripts that are executed sequentially. They cannot provide a dependency between scripts nor can they provide a restart of services when problems are corrected.
- **Oracle Solaris ZFS**—Oracle Solaris ZFS provides unparalleled data integrity, capacity, performance, and manageability for storage. ZFS provides high-resiliency features, such as metadata logging to guarantee data integrity and speed recovery in the event of system failure. What differentiates Oracle Solaris ZFS from other competitive file system offerings is strong data integrity and high resiliency.

Oracle Solaris ZFS dramatically simplifies file system administration to help increase protection against administrative error. Oracle Solaris ZFS uses techniques such as copy-on-write and end-to-end checksumming to keep data consistent and eliminate silent data corruption. Because the file system is always consistent, time-consuming recovery procedures such as using `fsck` are not required if the system is shut down in an unclean manner. In addition, data is read and checked constantly to help ensure correctness, and any errors detected in a mirrored pool are automatically repaired to protect against costly and time-consuming data loss and previously undetectable silent data corruption. Corrections are made possible by a RAID-Z implementation that uses parity, striping, and atomic operations to aid the reconstruction of corrupted data.

Oracle Solaris ZFS provides file system end-to-end data protection as well as facilities to self-correct data. It also offers a simplified management model that does not require a 3rd party volume manager. ZFS offers the scalability to grow to the largest data storage requirements. Oracle Solaris ZFS constantly reads and checks data to help ensure it is correct, and if it detects an error in a storage pool with redundancy (protected with mirroring, ZFS RAID-Z, or ZFS RAID-Z2), Oracle Solaris ZFS automatically repairs the corrupt data. It optimizes file system reliability by maintaining data redundancy on commodity hardware through the delivery of basic mirroring, compression, and integrated volume management.

- **Oracle Solaris DTrace**—DTrace, a feature of Oracle Solaris, is a dynamic tracing framework for troubleshooting systemic problems in real time on production systems. DTrace is designed to quickly identify the root cause of system performance problems. DTrace safely and dynamically instruments the running operating system (OS) kernel and running applications without rebooting the kernel and recompiling—or even restarting—applications. This design greatly improves service uptime.

The Predictive Self Healing capability of Oracle Solaris and the fault management combination of the Fault Management Architecture and the Service Management Facility can offline processor threads or

cores in faults, retire suspect pages of memory, log errors or faults from I/O, or deal with any other issue detected by the server.

## Oracle VM Server for SPARC

Supported in all servers from Oracle using Oracle's multicore/multithreaded technology, Oracle VM Server for SPARC provides full virtual machines that run an independent operating system instance. Each logical domain is completely isolated, and the maximum number of virtual machines created on a single platform relies upon the capabilities of the hypervisor, rather than on the number of physical hardware devices installed in the system.

Oracle VM Server for SPARC 3.0 has the ability to perform a live migration from one domain to another. As the term *live migration* implies, the source domain and application no longer need to be halted or stopped. Migration of a running application from one domain to another is now possible with Oracle VM Server for SPARC 3.0. This allows a logical domain on the server to be moved to a different Dynamic Domain on the same server, on a different SPARC M5-32 server, or to a server based on Oracle's SPARC T2, T3, T4, or T5 processors.

Logical domains can also host Oracle Solaris Zones to capture the isolation, flexibility, and manageability features of both technologies. Deeply integrating Oracle VM Server for SPARC with the SPARC M5 processor, Oracle Solaris increases flexibility, isolates workload processing, and improves the potential for maximum server utilization.

## Oracle Solaris Cluster

While SPARC M5-32 servers deliver high levels of availability, Oracle Solaris Cluster enables organizations to deliver highly available application services. To limit outages due to single points of failure, mission-critical services need to be run in clustered physical servers that efficiently and smoothly take over the services from failing nodes, with minimal interruption to data services. While SPARC M5-32 servers are designed with full redundancy at the hardware level, Oracle Solaris Cluster provides the best high availability (HA) solution for SPARC servers running Oracle Solaris and applications.

Oracle Solaris Cluster is focused on failover between zones and logical domains within the server as well as to external servers. Tightly coupled with Oracle Solaris, Oracle Solaris Cluster detects failures without delay (zero-second delay) and provides much faster failure notification, application failover, and reconfiguration time. Oracle Solaris Cluster offers HA for today's complex solution stacks, with failover protection from the application layer through to the storage layer. Figure 4 shows an example of the time it takes Oracle Solaris Cluster to detect a failure and recover on the redundant node. Compared to the competing vendor solution shown in the figure, Oracle Solaris Cluster offers 62 percent faster service recovery time.

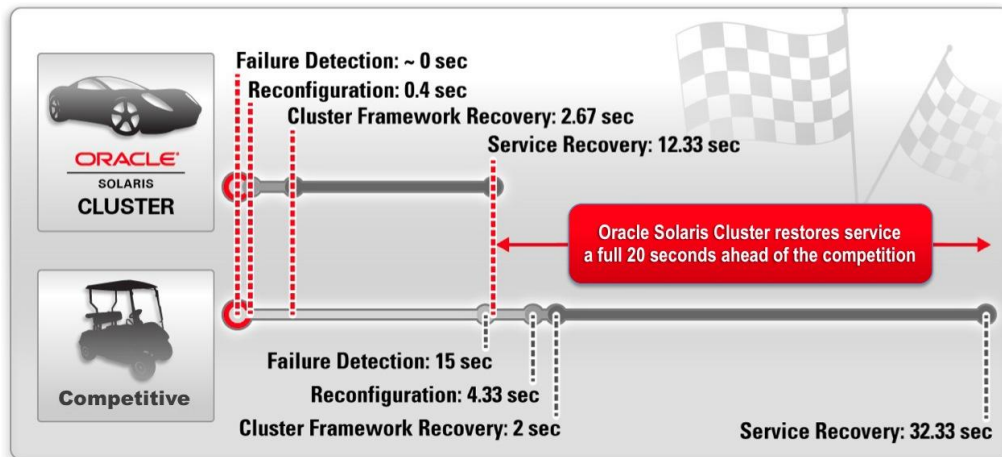


Figure 4. Oracle Solaris Cluster recovery times.

The use of Oracle Solaris Cluster enables much faster resumption of IT services. Oracle Solaris Cluster speeds this process on SPARC M5-32 servers and does the following:

- Integrates tightly with the Predictive Self Healing framework and supports applications controlled by the Service Management Facility in Oracle Solaris Zones and logical domains
- Makes extensive use of Oracle's storage management and volume management capabilities
- Supports Oracle Solaris ZFS as a failover file system and as a boot file system, allowing the use of ZFS as the single file system type used
- Leverages Oracle Solaris ZFS features such as pooled storage, built-in redundancy, and data integrity
- Integrates with Oracle Enterprise Manager Ops Center

## Integrated System Monitoring

The SPARC M5-32 server provides comprehensive monitoring and notifications to enable administrators to proactively detect and respond to problems with hardware and software components. With direct connectivity to the hardware components of the SPARC M5-32 server, Oracle Enterprise Manager Ops Center can alert administrators to hardware-related faults and log service requests automatically through integration with Oracle Auto Service Request for immediate review by Oracle Customer Support. Problems that would have required a combination of database, system, and storage administrators to detect in traditional systems can now be diagnosed in minutes because of integrated systems monitoring for the entire SPARC M5-32 server.

## Oracle Enterprise Manager Ops Center

Oracle Enterprise Manager Ops Center helps IT staff understand and manage every architectural layer—from bare metal to operating systems and applications. It provides a centralized interface for physical and virtual machine lifecycle management, from power-on to decommissioning. In addition, it offers IT administrators unique insight into the user experience, business transactions, and business services, helping administrators quickly detect changes in system health and troubleshoot issues across the entire environment.

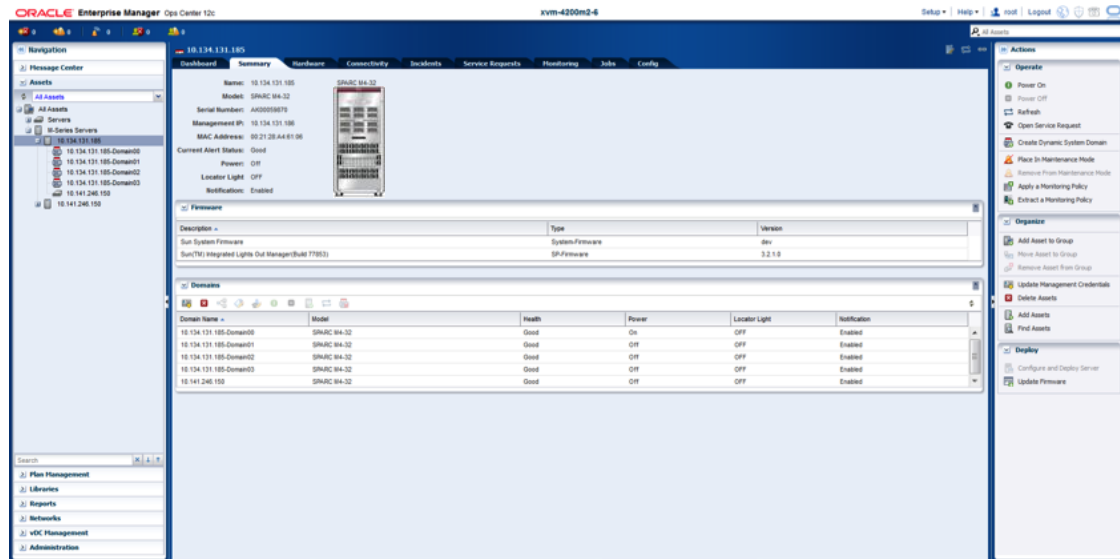


Figure 5. Oracle Enterprise Manager Ops Center user interface.

Oracle Enterprise Manager Ops Center is designed to discover and manage the SPARC M5-32 server as a whole system, not just as a collection of components in a rack. All hardware and software events are consolidated into Oracle Enterprise Manager Ops Center for a single system view of events. These events are then sent to My Oracle Support for further analysis and if needed, proactive action is taken. Oracle Enterprise Manager Ops Center delivers holistic fault management for events coming from the SPARC M5-32 server and even Oracle Solaris Cluster events are gathered.

## Oracle Auto Service Request

Once a problem is detected in a business-critical system, every minute until the problem is resolved is costly. Oracle Enterprise Manager Ops Center provides deep connections to My Oracle Support systems and processes to enable automatic problem detection, analysis, automated service requests, and access to the Oracle knowledge base and community for optimal problem resolution.

Oracle Auto Service Request is a feature of the Oracle hardware warranty and [Oracle Premier Support for Systems](#), which resolves problems by automatically opening service requests when specific hardware faults occur. Oracle Auto Service Request is integrated with [My Oracle Support](#). Customers must use My Oracle Support to activate Oracle Auto Service Request assets. Oracle Platinum Services

provides remote fault monitoring with rapid response times and patch deployment services to qualified Oracle Premier Support customers—at no additional cost.

## Conclusion

Oracle's SPARC M5-32 server is designed to help IT organizations consolidate multiple workloads in an environment that has been optimized for performance and RAS. It offers high levels of reliability and availability through redundant and highly reliable hardware components as well as multiple layers of software protection. The system strikes a perfect balance between the cost of hardware redundancy and the flexibility of software for data service availability.



Maximizing Reliability and Availability with the  
SPARC M5-32 Server  
March 2013, Version 1.0  
Author: Gary Combs

Oracle Corporation  
World Headquarters  
500 Oracle Parkway  
Redwood Shores, CA 94065  
U.S.A.

Worldwide Inquiries:  
Phone: +1.650.506.7000  
Fax: +1.650.506.7200

oracle.com



| Oracle is committed to developing practices and products that help protect the environment

Copyright © 2013, Oracle and/or its affiliates. All rights reserved. This document is provided for information purposes only and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. We specifically disclaim any liability with respect to this document and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission.

Oracle and Java are registered trademarks of Oracle and/or its affiliates. Other names may be trademarks of their respective owners.

Intel and Intel Xeon are trademarks or registered trademarks of Intel Corporation. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. AMD, Opteron, the AMD logo, and the AMD Opteron logo are trademarks or registered trademarks of Advanced Micro Devices. UNIX is a registered trademark of The Open Group. 1012

**Hardware and Software, Engineered to Work Together**