



Sun Cluster: Guía de conceptos para SO Solaris

Sun Microsystems, Inc.
4150 Network Circle
Santa Clara, CA 95054
U.S.A.

Referencia: 819-0162
September 2004, Revision A

Copyright 2004 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. Reservados todos los derechos.

Este producto o documento está protegido por la ley de copyright y se distribuye bajo licencias que restringen su uso, copia, distribución y descompilación. No se puede reproducir parte alguna de este producto o documento en ninguna forma ni por cualquier medio sin la autorización previa por escrito de Sun y sus licenciadores, si los hubiera. El software de otras empresas, incluida la tecnología de los tipos de letra, está protegido por la ley de copyright y con licencia de los distribuidores de Sun.

Determinadas partes del producto pueden derivarse de Berkeley BSD Systems, con licencia de la Universidad de California. UNIX es una marca registrada en los EE.UU. y otros países, bajo licencia exclusiva de X/Open Company, Ltd.

Sun, Sun Microsystems, el logotipo de Sun, docs.sun.com, AnswerBook, AnswerBook2, Sun Cluster, SunPlex, Sun Enterprise, Sun Enterprise 10000, Sun Enterprise SyMON, Sun Management Center, Solaris, Solaris Volume Manager, Sun StorEdge, Sun Fire, SPARCstation, OpenBoot y Solaris son marcas comerciales, marcas comerciales registradas o marcas de servicio de Sun Microsystems, Inc. en los EE.UU. y en otros países. Todas las marcas registradas SPARC se usan bajo licencia y son marcas comerciales o marcas registradas de SPARC International, Inc. en los EE.UU. y en otros países. Los productos con las marcas registradas de SPARC se basan en una arquitectura desarrollada por Sun Microsystems, Inc. ORACLE, Netscape

La interfaz gráfica de usuario OPEN LOOK y Sun™ fue desarrollada por Sun Microsystems, Inc. para sus usuarios y licenciatarios. Sun reconoce los esfuerzos pioneros de Xerox en la investigación y desarrollo del concepto de interfaces gráficas o visuales de usuario para la industria de la computación. Sun mantiene una licencia no exclusiva de Xerox para la interfaz gráfica de usuario de Xerox, que también cubre a los licenciatarios de Sun que implementen GUI de OPEN LOOK y que por otra parte cumplan con los acuerdos de licencia por escrito de Sun.

Derechos gubernamentales de los EE.UU. – Software comercial. Los usuarios del gobierno de los EE.UU. están sujetos a los acuerdos de la licencia estándar de Sun Microsystems Inc. y a las disposiciones aplicables sobre los FAR (derechos federales de adquisición) y sus suplementos.

ESTA DOCUMENTACIÓN SE PROPORCIONA "TAL CUAL". SE RENUNCIA A TODAS LAS CONDICIONES EXPRESAS O IMPLÍCITAS, REPRESENTACIONES Y GARANTÍAS, INCLUIDAS CUALQUIER GARANTÍA IMPLÍCITA DE COMERCIALIZACIÓN, ADECUACIÓN PARA UNA FINALIDAD DETERMINADA O DE NO CONTRAVENCIÓN, EXCEPTO EN AQUELLOS CASOS EN QUE DICHA RENUNCIA NO FUERA LEGALMENTE VÁLIDA.

Copyright 2004 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. Tous droits réservés.

Ce produit ou document est protégé par un copyright et distribué avec des licences qui en restreignent l'utilisation, la copie, la distribution, et la décompilation. Aucune partie de ce produit ou document ne peut être reproduite sous aucune forme, par quelque moyen que ce soit, sans l'autorisation préalable et écrite de Sun et de ses bailleurs de licence, s'il y en a. Le logiciel détenu par des tiers, et qui comprend la technologie relative aux polices de caractères, est protégé par un copyright et licencié par des fournisseurs de Sun.

Des parties de ce produit pourront être dérivées du système Berkeley BSD licenciés par l'Université de Californie. UNIX est une marque déposée aux Etats-Unis et dans d'autres pays et licenciée exclusivement par X/Open Company, Ltd.

Sun, Sun Microsystems, le logo Sun, docs.sun.com, AnswerBook, AnswerBook2, et Solaris sont des marques de fabrique ou des marques déposées, ou marques de service, de Sun Microsystems, Inc. aux Etats-Unis et dans d'autres pays. Toutes les marques SPARC sont utilisées sous licence et sont des marques de fabrique ou des marques déposées de SPARC International, Inc. aux Etats-Unis et dans d'autres pays. Les produits portant les marques SPARC sont basés sur une architecture développée par Sun Microsystems, Inc.

L'interface d'utilisation graphique OPEN LOOK et Sun™ a été développée par Sun Microsystems, Inc. pour ses utilisateurs et licenciés. Sun reconnaît les efforts de pionniers de Xerox pour la recherche et le développement du concept des interfaces d'utilisation visuelle ou graphique pour l'industrie de l'informatique. Sun détient une licence non exclusive de Xerox sur l'interface d'utilisation graphique Xerox, cette licence couvrant également les licenciés de Sun qui mettent en place l'interface d'utilisation graphique OPEN LOOK et qui en outre se conforment aux licences écrites de Sun.

CETTE PUBLICATION EST FOURNIE "EN L'ETAT" ET AUCUNE GARANTIE, EXPRESSE OU IMPLICITE, N'EST ACCORDEE, Y COMPRIS DES GARANTIES CONCERNANT LA VALEUR MARCHANDE, L'APTITUDE DE LA PUBLICATION A REPOUDRE A UNE UTILISATION PARTICULIERE, OU LE FAIT QU'ELLE NE SOIT PAS CONTREFAISANTE DE PRODUIT DE TIERS. CE DENI DE GARANTIE NE S'APPLIQUERAIT PAS, DANS LA MESURE OU IL SERAIT TENU JURIDIQUEMENT NUL ET NON AVENU.



041111@10082



Contenido

Prefacio	7
1 Introducción y visión general	13
Introducción al sistema SunPlex	14
Tres puntos de vista sobre el sistema SunPlex	15
Punto de vista del personal de instalación y reparación del hardware	15
Punto de vista del administrador de sistemas	16
Punto de vista del programador de aplicaciones	17
Tareas del sistema de SunPlex	19
2 Conceptos clave de los proveedores de servicio de hardware	21
Los componentes del hardware y el software de SunPlex	21
Nodos del clúster	22
Dispositivos multisistema	24
Discos locales	26
Medios extraíbles	26
Interconexión del clúster	27
Interfaces de red pública	27
Sistemas cliente	28
Dispositivos de acceso a la consola	28
Consola de administración	29
SPARC: Ejemplos de topología de Sun Cluster	29
SPARC: Topología de pares en clúster	30
SPARC: Topología par+n	31
SPARC: Topología n+1 (estrella)	32
SPARC: Topología n*n (escalable)	33

x86: Ejemplos de topología de Sun Cluster	34
x86: Topología de par en clúster	35
3 Conceptos clave de la administración y desarrollo de las aplicaciones	37
Interfaces administrativas	37
Hora del clúster	38
Estructura de alta disponibilidad (HA)	39
Supervisor de pertenencia al clúster	40
Depósito de configuración del clúster (CCR)	41
Dispositivos globales	41
ID de dispositivo (DID)	42
Grupos de dispositivos de discos	43
Recuperación de fallos del grupo de dispositivos de disco	43
Grupos de dispositivos de disco multipuerto	44
Espacio de nombres global	46
Ejemplo de espacios de nombres locales y globales	47
Sistemas de archivos del clúster	47
Uso de los sistemas de archivos del clúster	48
Tipo de recurso HAStoragePlus	49
Opción de montaje syncdir	49
Supervisión de las rutas de disco	50
Información general	50
Supervisión de las rutas del disco	52
Quórum y dispositivos de quórum	53
Acerca de los recuentos de votos de quórum	55
Aislamiento de fallos	55
Acerca de las configuraciones de quórum	57
Cumplimiento de los requisitos de dispositivos de quórum	58
Cumplimiento de las mejores prácticas recomendadas de dispositivos de quórum	58
Configuraciones de quórum recomendadas	60
Configuraciones de quórum atípicas	63
Configuraciones de quórum malas	64
Servicios de datos	65
Métodos de servicios de datos	67
Servicios de datos a prueba de fallos	68
Servicios de datos escalables	68
Valores de retroceso	71

Supervisores de fallos de servicios de datos	72
Desarrollo de nuevos servicios de datos	72
API de servicio de datos y de biblioteca de desarrollo de servicio de datos	73
Uso de la interconexión del clúster para el tráfico de servicio de datos	74
Recursos, grupos de recursos y tipos de recursos	76
Resource Group Manager (RGM)	77
Estados y configuración de recursos y grupos de recursos	77
Propiedades de recursos y grupos de recursos	79
Configuración del proyecto de servicios de datos	79
Determinación de requisitos para la configuración de proyectos	81
Establecimiento de límites de memoria virtual por proceso	82
Casos de recuperación de fallos	83
Adaptadores de red pública y IP Network Multipathing	89
SPARC: Compatibilidad con la reconfiguración dinámica	90
SPARC: Descripción general de la reconfiguración dinámica	91
SPARC: Consideraciones de la agrupación DR para dispositivos de CPU	91
SPARC: Consideraciones de la agrupación DR para memoria	92
SPARC: Consideraciones de la agrupación DR para unidades de disco y cinta	92
SPARC: Consideraciones de la agrupación DR para los dispositivos del quórum	92
SPARC: Consideraciones de la agrupación DR para interfaces de interconexión del clúster	93
SPARC: Consideraciones de la agrupación DR para interfaces de red pública	93
4 Preguntas más frecuentes (FAQ)	95
FAQ sobre alta disponibilidad	95
FAQ sobre sistemas de archivos	96
FAQ de gestión de volúmenes	97
FAQ de servicios de datos	98
FAQ de redes públicas	99
FAQ sobre pertenencia al clúster	100
FAQ de almacenamiento del clúster	101
FAQ de interconexión del clúster	101
FAQ sobre sistemas cliente	102
FAQ de consola de administración	102
FAQ sobre concentrador de terminal y procesador de servicio del sistema	103

Índice 105

Prefacio

Sun™ Cluster Concepts Guide for Solaris OS incluye información sobre conceptos y de referencia para sistemas SunPlex™ en sistemas basados en SPARC™ y x86.

Nota – En este documento el término x86'' hace referencia a la familia de chips microprocesadores Intel de 32 bits y a los compatibles de AMD.

El sistema SunPlex incluye todos los componentes de hardware y software que componen la solución clúster de Sun.

Este documento está pensado para administradores de sistema con experiencia que conozcan el software Sun Cluster, no se debe usar como guía de preventa ni de planificación. Antes de leerlo, debe conocer su sistema y disponer del equipo y el software adecuados.

Para entender los conceptos que se describen en este manual es necesario conocer el sistema operativo Solaris™ y tener experiencia con el software de gestión de volúmenes que utiliza el sistema SunPlex.

Nota – El software de Sun Cluster se ejecuta en dos plataformas: SPARC y x86. La información de este documento se aplica a ambas a menos que se especifique lo contrario de forma expresa en un capítulo, apartado, nota, elemento de lista, figura, tabla o ejemplo especiales.

Convenciones tipográficas

La tabla siguiente describe los cambios tipográficos utilizados en este manual.

TABLA P-1 Convenciones tipográficas

Tipo de letra o símbolo	Significado	Ejemplo
AaBbCc123	Los nombres de las órdenes, archivos, directorios y mensajes que aparecen en la pantalla del sistema	Modifique el archivo <code>.login</code> . Utilice el comando <code>ls -a</code> para mostrar todos los archivos. <code>nombre_sistema% tiene correo.</code>
AaBbCc123	Lo que usted escribe, contrastado con la salida por la pantalla del sistema	<code>nombre_máquina% su</code> Contraseña:
<i>AaBbCc123</i>	Plantilla de la línea de órdenes: sustituir por un valor o nombre real	El comando necesario para eliminar un archivo es <code>rm nombrearchivo</code> .
<i>AaBbCc123</i>	Títulos de los manuales, palabras y términos nuevos o palabras destacables	Consulte el capítulo 6 de la <i>Guía del usuario</i> . Éstas se denominan opciones de <i>clase</i> . <i>No</i> guarde el archivo. (En ocasiones se utiliza la negrita para enfatizar.)

Indicadores de los shells en los ejemplos de órdenes

La tabla siguiente muestra los indicadores predeterminados del sistema y de superusuario para los shells Bourne, Korn y C.

TABLA P-2 Indicadores de los shells

Shell	Indicador
Indicador del shell C	nombre_sistema%
Indicador de superusuario en el shell C	nombre_sistema#
Indicador de los shells Bourne y Korn	\$
Indicador de superusuario en los shell Bourne y Korn	#

Documentación relacionada

Puede encontrar información sobre temas referentes a Sun Cluster en la documentación enumerada en la tabla siguiente. Toda la documentación de Sun Cluster está disponible en <http://docs.sun.com>.

Tema	Documentación
Información general	<i>Sun Cluster Overview for Solaris OS</i>
Conceptos	<i>Sun Cluster Concepts Guide for Solaris OS</i>
Instalación y administración de hardware	<i>Sun Cluster 3.x Hardware Administration Manual for Solaris OS</i> Guías de administración de hardware individuales
Instalación del software	<i>Sun Cluster Software Installation Guide for Solaris OS</i>
Instalación y administración del servicio de datos	<i>Sun Cluster Data Services Planning and Administration Guide for Solaris OS</i> Guías de servicio de datos individuales
Desarrollo de los servicios de datos	<i>Sun Cluster Data Services Developer's Guide for Solaris OS</i>
Administración de sistema	<i>Sun Cluster System Administration Guide for Solaris OS</i>
Mensajes de error	<i>Sun Cluster Error Messages Guide for Solaris OS</i>
Referencias sobre las órdenes y las funciones	<i>Sun Cluster Reference Manual for Solaris OS</i>

Si desea una lista completa de la documentación sobre Sun Cluster, consulte las notas sobre la versión de Sun Cluster en <http://docs.sun.com>.

Acceso a la documentación de Sun en línea

La sede web docs.sun.comSM permite acceder a la documentación técnica de Sun en línea. Puede explorar el archivo docs.sun.com, buscar el título de un manual o un tema específicos. El URL es <http://docs.sun.com>.

Solicitud de documentación de Sun

Sun Microsystems ofrece una seleccionada documentación impresa sobre el producto. Para obtener una lista de documentos y la forma como adquirirlos, consulte "Adquirir documentación impresa" en <http://docs.sun.com>.

Obtención de ayuda

Si tiene problemas durante la instalación o utilización del sistema SunPlex, póngase en contacto con su proveedor de servicios y proporcione la información siguiente:

- Su nombre y dirección de correo electrónico (si estuviera disponible)
- El nombre, dirección y número de teléfono de su empresa
- Los modelos y números de serie de sus sistemas
- El número de versión del sistema operativo; por ejemplo Solaris 9
- El número de versión del software Sun Cluster (por ejemplo, 3.1 4/04)

Use las órdenes siguientes para reunir información sobre todos los nodos del sistema para el proveedor de servicio.

Comando	Función
<code>prtconf -v</code>	Muestra el tamaño de la memoria del sistema y ofrece información sobre los dispositivos periféricos
<code>psrinfo -v</code>	Muestra información acerca de los procesadores

Comando	Función
<code>showrev -p</code>	Indica las modificaciones instaladas
<code>SPARC: prtdiag -v</code>	Muestra información de diagnóstico del sistema
<code>scinstall -pv</code>	Muestra información sobre la versión y el paquete de Sun Cluster.
<code>scstat</code>	Ofrece una instantánea del estado del clúster.
<code>scconf -p</code>	Muestra información de configuración sobre el clúster
<code>scrgadm -p</code>	Muestra información sobre los recursos instalados y los grupos y tipos de recursos

Tenga también a punto el contenido del archivo `/var/adm/messages`.

Introducción y visión general

El sistema SunPlex es una solución integrada de hardware y software Sun Cluster que se utiliza para crear servicios de alta disponibilidad y escalabilidad.

El manual Sun Cluster: Guía de conceptos para SO Solaris proporciona la información conceptual que necesitan los usuarios principales de la documentación de SunPlex, entre los que se incluyen:

- Proveedores de servicio que instalan y reparan hardware de clúster
- Administradores de sistemas que instalan, configuran y administran software de Sun Cluster
- Desarrolladores de aplicaciones que desarrollan servicios de recuperación de fallos y escalables para aplicaciones no incluidas actualmente en el producto Sun Cluster

Este manual complementa el resto de la documentación de SunPlex para proporcionar una visión completa del sistema SunPlex.

Este capítulo

- Proporciona una introducción y una visión general del sistema SunPlex
- Describe los distintos puntos de vista de la audiencia de SunPlex
- Identifica conceptos clave que es necesario entender antes de trabajar con el sistema SunPlex
- Correlaciona los conceptos clave con la documentación de SunPlex que incluye procedimientos e información relacionada
- Correlaciona tareas relacionadas del clúster con la documentación que contiene procedimientos usados para realizar esas tareas

Introducción al sistema SunPlex

El sistema SunPlex amplía el sistema operativo Solaris hasta hacerlo un sistema operativo de clúster. Un clúster, o plex, es una colección de nodos informáticos acoplados indirectamente que proporcionan una vista de cliente único de servicios de red o de aplicaciones, incluidos bases de datos, servicios de web y servicios de archivos.

Cada nodo del clúster es un servidor autónomo que ejecuta sus propios procesos. Éstos se comunican entre sí para formar lo que parece (a un cliente de red) como un sólo sistema que proporciona de forma cooperativa aplicaciones, recursos de sistema y datos a usuarios.

Un clúster ofrece, respecto a los sistemas tradicionales de servidor único, varias ventajas que incluyen soporte para servicios de protección contra fallos y escalabilidad, capacidad para crecimiento modular y un precio básico económico en comparación a los sistemas de hardware tolerantes a fallos.

Los objetivos del sistema SunPlex son:

- Reducir o eliminar el tiempo de inactividad del sistema debido a los fallos de software o hardware
- Asegurar la disponibilidad de los datos y aplicaciones a los usuarios finales, cualquiera que sea el tipo de fallo que normalmente dejaría inactivo un sistema de servidor único
- Aumentar el rendimiento de las aplicaciones permitiendo escalar procesadores adicionales con solo agregar más nodos al clúster
- Proporcionar una mejor disponibilidad del sistema al permitirle realizar el mantenimiento sin tener que apagar todo el clúster

Para obtener información sobre la tolerancia a los fallos y la alta disponibilidad, consulte “Making Applications Highly Available With Sun Cluster” en *Sun Cluster Overview for Solaris OS*.

Consulte “FAQ sobre alta disponibilidad” en la página 95 para consultar preguntas y respuestas referentes a la alta disponibilidad.

Tres puntos de vista sobre el sistema SunPlex

Este apartado describe tres puntos de vista distintos sobre el sistema SunPlex, los conceptos clave y la documentación relevante para cada punto de vista. Estos puntos de vista provienen de:

- Personal de instalación y reparación de hardware
- Administradores del sistema
- Programadores de aplicaciones

Punto de vista del personal de instalación y reparación del hardware

Para los profesionales de la reparación del hardware, el sistema SunPlex es como una colección de hardware estándar que incluye servidores, redes y almacenamiento. Estos componentes están todos interconectados de manera que todos tengan un recambio y no exista un punto débil para los fallos.

Conceptos clave del hardware

El personal de reparación de hardware necesita entender los conceptos de clúster siguientes.

- Configuraciones de hardware de clúster y cableado
- Instalación y reparación (agregar, eliminar, sustituir):
 - Componentes de la interfaz de red (adaptadores, uniones, cables)
 - Tarjetas interfaz de disco
 - Matrices de disco
 - Unidades de disco
 - La consola de administración y el dispositivo de acceso a la consola
- Configuración de la consola de administración y del dispositivo de acceso a la consola

Referencias conceptuales de hardware recomendadas

Los apartados siguientes contienen material relacionado con los conceptos clave anteriores:

- [“Nodos del clúster” en la página 22](#)

- “Dispositivos multisistema” en la página 24
- “Discos locales” en la página 26
- “Interconexión del clúster” en la página 27
- “Interfaces de red pública” en la página 27
- “Sistemas cliente” en la página 28
- “Consola de administración” en la página 29
- “Dispositivos de acceso a la consola” en la página 28
- “SPARC: Topología de pares en clúster” en la página 30
- “SPARC: Topología n+1 (estrella)” en la página 32

Información destacable sobre SunPlex

El documento de SunPlex siguiente incluye procedimientos e información asociada con conceptos de servicios de hardware:

Sun Cluster 3.x Hardware Administration Manual for Solaris OS

Punto de vista del administrador de sistemas

Para el administrador, el sistema SunPlex se asemeja a un conjunto de servidores (nodos) interconectados y que comparten dispositivos de almacenamiento. El administrador del sistema ve:

- Software para clústers especializado que se integra con el de Solaris para supervisar la conectividad entre los nodos del clúster
- Software especializado que supervisa el buen funcionamiento de los programas de aplicación del usuario que se ejecutan en los nodos del clúster
- Software de gestión de volúmenes que configura y administra discos
- Software de clúster especializado que permite a todos los nodos acceder a todos los dispositivos de almacenamiento, incluso a aquellos no conectados directamente a discos
- Software de clúster especializado que permite a los archivos aparecer en todos los nodos como si estuvieran alojados localmente en él

Conceptos clave de la administración de sistemas

Es importante que los administradores de sistemas comprendan los conceptos y procesos siguientes:

- La interacción entre los componentes de hardware y software
- El flujo general sobre cómo instalar y configurar el clúster, que incluye:
 - Instalación del sistema operativo Solaris

- Instalación y configuración del software Sun Cluster
- Instalación y configuración de un gestor de volúmenes
- Instalación y configuración del software de aplicación para que esté preparado para el clúster
- Instalación y configuración del software del servicio de datos de Sun Cluster
- Procedimientos administrativos del clúster para agregar, eliminar, sustituir y reparar componentes de hardware y software del clúster
- Modificaciones de la configuración para mejorar el rendimiento

Referencias conceptuales recomendadas para administradores de sistemas

Los apartados siguientes contienen material relacionado con los conceptos clave anteriores:

- “Interfaces administrativas” en la página 37
- “Hora del clúster” en la página 38
- “Estructura de alta disponibilidad (HA)” en la página 39
- “Dispositivos globales” en la página 41
- “Grupos de dispositivos de discos” en la página 43
- “Espacio de nombres global” en la página 46
- “Sistemas de archivos del clúster” en la página 47
- “Supervisión de las rutas de disco” en la página 50
- “Aislamiento de fallos” en la página 55
- “Servicios de datos” en la página 65

Documentación de SunPlex importante para el administrador de sistemas

Los documentos de SunPlex siguientes incluyen procedimientos e información asociada a los conceptos de la administración de sistemas:

- *Sun Cluster Software Installation Guide for Solaris OS*
- *Sun Cluster System Administration Guide for Solaris OS*
- *Sun Cluster Error Messages Guide for Solaris OS*
- *Sun Cluster 3.1 9/04 Release Notes for Solaris OS*
- *Sun Cluster 3.x Release Notes Supplement*

Punto de vista del programador de aplicaciones

El sistema SunPlex proporciona *servicios de datos* para aplicaciones como Oracle (en sistemas basados en plataformas SPARC), NFS, DNS, Sun™ Java System Web Server (anteriormente Sun Java System Web Server), Apache Web Server (en sistemas basados en plataformas SPARC) y Sun Java System Directory Server (anteriormente Sun Java

System Directory Server). Los servicios de datos se crean configurando aplicaciones estándar de forma que se ejecuten bajo el control del software Sun Cluster. Éste proporciona archivos de configuración y métodos de gestión que inician, paran y supervisan las aplicaciones. Si necesita crear un nuevo servicio a prueba de fallos o escalable, puede usar la Interfaz de programación de aplicaciones (API) de SunPlex y la API de tecnologías de habilitación de servicio de datos (DSET API) para desarrollar los archivos de configuración y métodos de gestión necesarios que permitan a esa aplicación ejecutarse como servicio de datos en el clúster.

Conceptos clave del programador de aplicaciones

Es importante que los programadores de aplicaciones entiendan lo siguiente:

- Las características de la aplicación, para determinar si ésta puede ejecutarse como servicio de recuperación de fallos o de escalabilidad de datos.
- La API de Sun Cluster, API DSET y el servicio de datos “genérico”. Los programadores deben determinar qué herramientas son más adecuadas para usarlas al escribir programas o secuencias que configuren su aplicación para el entorno clúster.

Referencias conceptuales recomendadas para los programadores de aplicaciones

Los apartados siguientes contienen material relacionado con los conceptos clave anteriores:

- [“Servicios de datos” en la página 65](#)
- [“Recursos, grupos de recursos y tipos de recursos” en la página 76](#)
- [Capítulo 4](#)

Documentación sobre SunPlex importante para el programador de aplicaciones

Los documentos de SunPlex siguientes incluyen procedimientos e información asociada a los conceptos de la programación de aplicaciones:

- *Sun Cluster Data Services Developer’s Guide for Solaris OS*
- *Sun Cluster Data Services Planning and Administration Guide for Solaris OS*

Tareas del sistema de SunPlex

Todas las tareas del sistema SunPlex requieren algunos conocimientos generales previos. La tabla siguiente incluye una visión general de las tareas y la documentación que describe los pasos de las tareas. El apartado de conceptos de este manual describe la correlación entre los conceptos y estas tareas.

TABLA 1-1 Mapa de tareas: correspondencia entre las tareas del usuario y la documentación

Para realizar esta tarea...	Usar esta documentación...
Instalar el hardware del clúster	<i>Sun Cluster 3.x Hardware Administration Manual for Solaris OS</i>
Instalar el software de Solaris en el clúster	<i>Sun Cluster Software Installation Guide for Solaris OS</i>
SPARC: Instalar el software Sun™ Management Center	<i>Sun Cluster Software Installation Guide for Solaris OS</i>
Instalar y configurar el software Sun Cluster	<i>Sun Cluster Software Installation Guide for Solaris OS</i>
Instalar y configurar el software de gestión de volúmenes	<i>Sun Cluster Software Installation Guide for Solaris OS</i> La documentación del gestor de volúmenes
Instalar y configurar los servicios de datos para Sun Cluster	<i>Sun Cluster Data Services Planning and Administration Guide for Solaris OS</i>
Reparar el hardware del clúster	<i>Sun Cluster 3.x Hardware Administration Manual for Solaris OS</i>
Administrar el software de Sun Cluster	<i>Sun Cluster System Administration Guide for Solaris OS</i>
Administrar el software de la gestión de volúmenes	<i>Sun Cluster System Administration Guide for Solaris OS</i> y la documentación de gestión de volúmenes
Administrar el software de las aplicaciones	La documentación de su aplicación
Identificar los problemas y las acciones de usuario sugeridas	<i>Sun Cluster Error Messages Guide for Solaris OS</i>
Crear un servicio de datos nuevo	<i>Sun Cluster Data Services Developer's Guide for Solaris OS</i>

Conceptos clave de los proveedores de servicio de hardware

Este capítulo describe los conceptos clave relacionados con los componentes de hardware de una configuración de sistema de SunPlex. Se tratan los siguientes temas:

- “Nodos del clúster” en la página 22
- “Dispositivos multisistema” en la página 24
- “Discos locales” en la página 26
- “Medios extraíbles” en la página 26
- “Interconexión del clúster” en la página 27
- “Interfaces de red pública” en la página 27
- “Sistemas cliente” en la página 28
- “Dispositivos de acceso a la consola” en la página 28
- “Consola de administración” en la página 29
- “SPARC: Ejemplos de topología de Sun Cluster” en la página 29
- “x86: Ejemplos de topología de Sun Cluster” en la página 34

Los componentes del hardware y el software de SunPlex

Esta información está dirigida principalmente a los proveedores de servicio de hardware a quienes ayuda a entender la relación entre los componentes de hardware para mejor instalar, configurar o reparar el hardware del clúster. Asimismo, a los administradores del sistema del clúster esta información les puede resultar útil para instalar, configurar y administrar el software del clúster.

Un clúster consta de varios componentes de hardware, entre los que se incluyen:

- Nodos de clúster con discos locales (no compartidos)
- Almacenamiento multisistema (discos compartidos entre nodos)
- Medios extraíbles (cintas y CD-ROM)

- Interconexión del clúster
- Interfaces de red públicas
- Sistemas cliente
- Consola de administración
- Dispositivos de acceso a la consola

El sistema SunPlex permite combinar estos componentes en diversas configuraciones, que se describen en [“SPARC: Ejemplos de topología de Sun Cluster”](#) en la página 29.

Para ver una ilustración de un ejemplo de configuración de un clúster de dos nodos, consulte [“Sun Cluster Hardware Environment”](#) en *Sun Cluster Overview for Solaris OS*.

Nodos del clúster

Un nodo del clúster es una máquina que ejecuta el sistema operativo Solaris y el software de Sun Cluster y es un componente actual del clúster (un *miembro del clúster*) o un miembro potencial.

SPARC: El software de Sun Cluster permite tener de dos a ocho nodos por clúster. Consulte en [“SPARC: Ejemplos de topología de Sun Cluster”](#) en la página 29 las configuraciones de nodos admitidas.

x86: Sun Cluster permite tener dos nodos en un clúster. Consulte en [“x86: Ejemplos de topología de Sun Cluster”](#) en la página 34 las configuraciones de nodos admitidas.

Los nodos del clúster normalmente están conectados a uno o más dispositivos multisistema, aquéllos que no lo están usan el sistema de archivos del clúster para acceder a los dispositivos multisistema. Por ejemplo, una configuración de servicio escalable permite que los nodos atiendan peticiones sin estar directamente conectados a los dispositivos multisistema.

Además, los nodos en configuraciones de bases de datos paralelas comparten acceso simultáneo a todos los discos. Consulte [“Dispositivos multisistema”](#) en la página 24 y el [Capítulo 3](#) para obtener más información sobre las configuraciones de bases de datos paralelas.

Todos los nodos del clúster están agrupados bajo un nombre común (el nombre del clúster) que se utiliza para acceder a éste y gestionarlo.

Los adaptadores de redes públicas conectan los nodos a éstas para ofrecer acceso de cliente al clúster.

Los miembros del clúster se comunican entre sí a través de una o más redes físicamente independientes que reciben el nombre de *interconexión del clúster*.

Todos los nodos del clúster reciben información de cuándo un nodo se une o deja el clúster, conocen también los recursos que se están ejecutando localmente así como los que se están ejecutando en los otros nodos del clúster.

Los nodos del mismo clúster deben tener capacidades de procesamiento, memoria y E/S similares para permitir que la recuperación de fallos se produzca sin que haya una degradación importante en el rendimiento. Debido a la posibilidad de recuperación de fallos, todos los nodos deben tener suficiente capacidad sobrante para que los que sean de respaldo o secundarios puedan aprovechar la carga de trabajo.

Cada nodo ejecuta su propio sistema de archivos raíz individual (/).

Componentes del software para los miembros del hardware del clúster

Para poder actuar como miembro del clúster, es necesario tener instalado el software siguiente:

- Sistema operativo Solaris
- Software Sun Cluster
- Aplicación de servicio de datos
- Gestión de volúmenes (Solaris Volume Manager™ o VERITAS Volume Manager)
Una excepción es la configuración que usa una matriz redundante de hardware de discos independientes (RAID). Esta configuración puede no requerir un software de gestión de volúmenes del tipo Solaris Volume Manager o VERITAS Volume Manager.
- Consulte *Sun Cluster Software Installation Guide for Solaris OS* para obtener más información sobre cómo instalar el sistema operativo Solaris, Sun Cluster, y el software de gestión de volúmenes.
- Consulte *Sun Cluster Data Services Planning and Administration Guide for Solaris OS* para obtener información sobre cómo instalar y configurar los servicios de datos.
- Consulte el [Capítulo 3](#) para obtener información básica sobre los componentes del software mencionados anteriormente.

La figura siguiente incluye una vista detallada de los componentes del software que funcionan juntos para crear el entorno Sun Cluster.

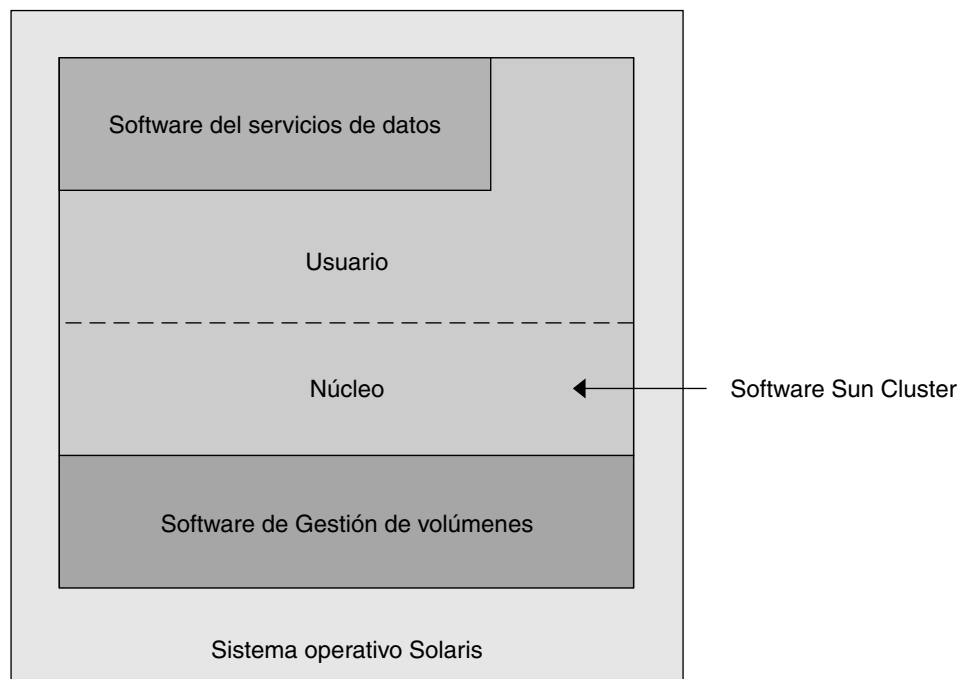


FIGURA 2-1 Relaciones entre los componentes del software de Sun Cluster

Consulte el [Capítulo 4](#) para obtener preguntas y respuestas sobre los miembros de clúster.

Dispositivos multisistema

Los discos que pueden conectarse a más de un nodo simultáneamente se denominan multisistema. En el entorno Sun Cluster, el almacenamiento multisistema hace que los discos estén muy disponibles. Sun Cluster requiere que haya almacenamiento multisistema en dos clústers del nodo para establecer el quórum. Los clústers de más de tres nodos no requieren almacenamiento multisistema.

Los dispositivos multisistema poseen las características siguientes.

- Pueden resistir los fallos de nodos individuales.
- Almacenan datos de las aplicaciones y también pueden almacenar los archivos binarios y de configuración de éstas.
- Protegen contra los fallos del nodo. Si las peticiones de los clientes están accediendo a datos a través de un nodo y éste falla, las peticiones se desvían para usar otro nodo que tenga una conexión directa con los mismos discos.

- A ellos se accede tanto globalmente a través de un nodo primario que “coordina” los discos como por acceso directo simultáneo a través de las rutas de acceso locales. Actualmente la única aplicación que usa acceso directo simultáneo es Oracle Real Application Clusters.

Un gestor de volúmenes proporciona configuraciones especulares o RAID-5 para la redundancia de los dispositivos multisistema. Actualmente, Sun Cluster admite Solaris Volume Manager™ y VERITAS Volume Manager, disponibles sólo para su uso en clústers basados en plataformas SPARC, como gestores de volúmenes, y el controlador de hardware RDAC RAID-5 en algunas plataformas de hardware RAID.

Al combinar los dispositivos multisistema con los discos especulares y la separación de datos se consigue proteger contra el fallo de los nodos y el de los discos individuales.

Consulte el [Capítulo 4](#) para obtener preguntas y respuestas sobre el almacenamiento multisistema.

Iniciador múltiple SCSI

Este apartado se aplica sólo a dispositivos de almacenamiento SCSI y no a los de fibra óptica que usan los dispositivos multisistema.

En un servidor autónomo, el nodo servidor controla las actividades del bus SCSI porque el circuito adaptador del sistema SCSI se conecta a ese servidor en un bus SCSI determinado. Al circuito adaptador del sistema SCSI se le conoce como *iniciador SCSI*. Este circuito inicia todas las actividades de este bus SCSI. La dirección SCSI predeterminada de los adaptadores en el sistema Sun es 7.

Las configuraciones del clúster comparten el almacenamiento entre varios nodos del servidor, usando los dispositivos multisistema. Cuando el almacenamiento del clúster se compone de dispositivos SCSI de terminación simple o diferencial a la configuración se la denomina iniciador múltiple SCSI. Tal como implica esta terminología, existe más de un iniciador SCSI en el bus.

La especificación SCSI requiere que cada dispositivo de un bus SCSI tenga una dirección exclusiva. (El adaptador del sistema también es un dispositivo del bus SCSI.) La configuración predeterminada del hardware en un entorno de iniciador múltiple provoca un conflicto debido a que de forma predeterminada todos los adaptadores del sistema SCSI tienen el valor 7.

Para resolver este conflicto, se debe dejar uno de los adaptadores de cada bus SCSI con la dirección SCSI 7 y se debe configurar los otros adaptadores como direcciones SCSI no utilizadas. La planificación correcta dictamina que estas direcciones SCSI “no utilizadas” incluyan tanto direcciones usadas actualmente como no utilizadas. Un ejemplo de direcciones no utilizadas en el futuro es la incorporación del almacenamiento instalando unidades nuevas en ranuras de unidad vacías.

En la mayoría de configuraciones, la dirección SCSI disponible para un segundo adaptador de sistema es 6.

Se pueden cambiar las direcciones SCSI seleccionadas de esos adaptadores del sistema mediante una de las siguientes herramientas de configuración de la propiedad `scsi-initiator-id`:

- `eeeprom(1M)`
- PROM de OpenBoot en un sistema basado en la plataforma SPARC
- La utilidad SCSI que se puede ejecutar opcionalmente después de que arranque la BIOS en un sistema basado en la plataforma x86

Se puede configurar esta propiedad tanto globalmente para un nodo como individualmente para cada adaptador. En el capítulo para cada alojamiento de disco de *Sun Cluster Hardware Collection* hay instrucciones para configurar un `scsi-initiator-id` exclusivo para cada adaptador SCSI.

Discos locales

Los discos locales son aquellos que sólo están conectados a un nodo individual. Por tanto, no están protegidos contra ningún fallo del nodo (no son de alta disponibilidad). A pesar de ello, todos los discos, incluso los locales se incluyen en el espacio de nombres global y se configuran como *dispositivos globales*. Por lo tanto, los discos en sí son visibles desde todos los nodos del clúster.

Puede poner los sistemas de archivos de discos locales a disposición de otros nodos situándolos bajo un punto de montaje global. Si el nodo que tiene montado actualmente uno de estos sistemas de archivos globales falla, el resto de nodos pierde acceso a ese sistema de archivos. Usar un gestor de volúmenes permite duplicar estos discos de forma que un fallo no pueda provocar que los sistemas de archivos dejen de estar accesibles, aunque los gestores de volúmenes no protejan contra los fallos de los nodos.

Consulte el apartado “[Dispositivos globales](#)” en la [página 41](#) para obtener más información sobre dispositivos globales.

Medios extraíbles

El clúster admite medios extraíbles, como unidades de cinta y de CD-ROM, que se instalan, configuran y reparan de la misma forma que en un entorno que no está configurado en clúster. Estos dispositivos están configurados como globales en Sun Cluster, de manera que son accesibles desde cualquier nodo del clúster. Consulte *Sun Cluster 3.x Hardware Administration Manual for Solaris OS* para obtener información sobre la instalación y configuración de medios extraíbles.

Consulte el apartado “[Dispositivos globales](#)” en la [página 41](#) para obtener más información sobre dispositivos globales.

Interconexión del clúster

La *interconexión del clúster* es la configuración física de dispositivos que se utiliza para transferir comunicaciones privadas del clúster y de servicios de datos entre los nodos del clúster. Debido a que la interconexión se utiliza ampliamente para comunicaciones privadas del clúster, puede limitar el rendimiento.

La interconexión del clúster sólo puede conectar los distintos nodos del clúster. El modelo de seguridad Sun Cluster asume que sólo los nodos del clúster tienen acceso físico a la interconexión del clúster.

Todos los nodos deben estar conectados por la interconexión del clúster a través de al menos dos redes independientes físicamente o rutas de acceso, para evitar que exista un único punto de fallo. Entre dos nodos cualesquiera puede tener varias redes independientes físicamente (de dos a seis). La interconexión del clúster se compone de tres componentes del hardware: adaptadores, uniones y cables.

La lista siguiente describe cada uno de estos componentes del hardware.

- **Adaptadores:** las tarjetas de interfaz de red que residen en cada nodo del clúster. Sus nombres se construyen a partir de un nombre de dispositivo seguido inmediatamente por un número de unidad física, por ejemplo: qfe2. Algunos adaptadores sólo tienen una conexión de red física, aunque otros como la tarjeta qfe tienen varias conexiones físicas. Algunos también contienen interfaces de red y de almacenamiento.

Un adaptador de red con varias interfaces podría convertirse en un único punto de fallo si falla todo el adaptador. Para una máxima disponibilidad, se debe planificar el clúster de manera que la única ruta entre dos nodos no dependa de un único adaptador de red individual.

- **Uniones:** los conmutadores que residen fuera de los nodos del clúster. Llevan a cabo funciones de transporte y conmutación para permitir conectar más de dos nodos simultáneamente. En un clúster de dos nodos no hacen falta uniones porque los nodos pueden conectarse directamente entre sí con cables físicos redundantes conectados a adaptadores redundantes de cada nodo. Las configuraciones superiores a dos nodos en general requieren uniones.
- **Cables:** las conexiones físicas que van entre dos adaptadores de red o entre un adaptador y una unión.

Consulte el [Capítulo 4](#) para obtener preguntas y respuestas sobre la interconexión del clúster.

Interfaces de red pública

Los clientes se conectan al clúster a través de interfaces de red pública. Todas las tarjetas adaptadoras de red pueden conectarse a una o más redes públicas, de acuerdo con las interfaces de hardware que la tarjeta tenga. Los nodos pueden configurarse

para que incluyan múltiples tarjetas interfaz de red pública a fin de que varias estén activas y se utilicen en caso de recuperación de fallos como respaldo entre ellas. Si uno de los adaptadores falla, se llama al software IP Network Multipathing para la recuperación de la interfaz defectuosa pasando a otro adaptador del grupo.

No hay consideraciones de hardware especiales relacionadas con la agrupación en clúster para las interfaces de red públicas.

Consulte el [Capítulo 4](#) para obtener preguntas y respuestas sobre las redes públicas.

Sistemas cliente

Los sistemas cliente son las estaciones de trabajo u otros servidores que acceden al clúster a través de la red pública. Los programas del lado del cliente usan datos u otros servicios proporcionados por aplicaciones del lado del servidor que se ejecutan en el clúster.

Los sistemas cliente no son de alta disponibilidad. Los datos y aplicaciones del clúster son de alta disponibilidad.

Consulte el [Capítulo 4](#) para obtener preguntas y respuestas sobre los sistemas cliente.

Dispositivos de acceso a la consola

Es necesario que disponga de acceso a la consola para todos los nodos del clúster. Para ello, use el concentrador de terminal que ha adquirido con el hardware del clúster, el procesador de servicio del sistema (SSP) de los servidores Sun Enterprise E10000™, el controlador del sistema en servidores Sun Fire™ u otro dispositivo que pueda acceder a `ttys` en todos los nodos.

Sun sólo dispone de un concentrador de terminal admitido y su uso es opcional. El concentrador de terminal permite el acceso a `/dev/console` en todos los nodos a través de una red TCP/IP. El resultado es el acceso a nivel de consola para todos los nodos desde una estación de trabajo remota desde cualquier lugar de la red.

El procesador de servicio del sistema (SSP) proporciona acceso a la consola para el Sun Enterprise E10000 server. SSP es una máquina de una red Ethernet que está configurada para admitir el Sun Enterprise E10000 server. SSP es la consola de administración para el Sun Enterprise E10000 server. Gracias a la función Sun Enterprise E10000 Network Console, cualquier estación de trabajo de la red puede abrir una sesión de consola de sistema.

Otros métodos de acceso a la consola son concentradores de otras terminales, acceso de puerto serie por `tip` (1) desde otro nodo y terminales no inteligentes. Puede usar teclados y monitores de Sun™ u otros dispositivos de puerto serie si el proveedor de servicio de hardware los admite.

Consola de administración

Para administrar el clúster activo se puede emplear una estación de trabajo UltraSPARC® exclusiva o un servidor Sun Fire™ V65x, conocidos como *consola de administración*. Normalmente en la consola de administración se instala y se ejecuta el software de herramientas de administración, como el panel de control del clúster (CCP) y el módulo de Sun Cluster para el producto Sun Management Center™ (para su uso en clústers basados en la plataforma SPARC solamente). Si utiliza `cconsole` en CCP se le permitirá conectar a más de un nodo simultáneamente. Para obtener más información sobre el uso de CCP, consulte *Sun Cluster System Administration Guide*.

La consola de administración normalmente no es un nodo del clúster; se acostumbra a usar para obtener acceso remoto a los nodos del clúster, bien a través una red pública bien, opcionalmente, a través de un concentrador de terminales situado en una red. Si el clúster se compone de la plataforma Sun Enterprise E10000, es necesario que se pueda iniciar la sesión desde la consola de administración con el procesador de servicio del sistema (SSP) y que se pueda conectar mediante la orden `netcon (1M)`.

Normalmente los nodos se configuran sin monitor. A continuación se accede a la consola del nodo empleando una sesión `telnet` desde la consola de administración, que está conectada a un concentrador de terminal y desde éste al puerto serie del nodo. (En el caso del Sun Enterprise E10000 server, se conecta desde el procesador de servicio del sistema.) Para obtener más información, véase “Dispositivos de acceso a la consola” en la página 28.

Sun Cluster no requiere una consola de administración exclusiva, aunque si se usa una se obtendrán las siguientes ventajas:

- Permite la gestión centralizada del clúster ya que agrupa herramientas de consola y gestión en la misma máquina
- Ofrece al proveedor de servicio de hardware una resolución de problemas potencialmente más rápida

Consulte el [Capítulo 4](#) para obtener preguntas y respuestas sobre la consola de administración.

SPARC: Ejemplos de topología de Sun Cluster

Una topología es el esquema de conexión que une los nodos del clúster con las plataformas de almacenamiento que se usan en éste. Sun Cluster admite cualquier topología que cumpla las siguientes directrices.

- Sun Cluster, cuando se compone de sistemas basados en plataformas SPARC, admite un máximo de ocho nodos por clúster, sin tener en cuenta las configuraciones de almacenamiento que se hayan implementado.
- Un dispositivo de almacenamiento compartido puede conectarse a tantos nodos como se admitan.
- Los dispositivos de almacenamiento compartido no necesitan conectarse con todos los nodos del clúster. Sin embargo, estos dispositivos de almacenamiento deben conectarse a un mínimo de dos nodos.

Sun Cluster no obliga a configurar el clúster mediante topologías específicas. Las topologías que se describen a continuación se incluyen para proporcionar el vocabulario que permita discutir un esquema de conexión del clúster. Estas topologías son esquemas de conexión típicos.

- Pares en clúster
- Par+N
- N+1 (estrella)
- N*N (escalable)

Los apartados siguientes incluyen diagramas de ejemplo para cada topología.

SPARC: Topología de pares en clúster

Una topología de pares en clúster son dos o más pares de nodos que funcionan bajo un único marco administrativo de clúster. En esta configuración sólo se produce recuperación de fallos entre pares. Sin embargo, todos los nodos están conectados por la interconexión del clúster y funcionan bajo el control del software Sun Cluster. Esta topología se podría usar para ejecutar una aplicación de base de datos paralela en un par y una aplicación de recuperación de fallos o de escalabilidad en otro par.

Con el sistema de archivos del clúster, también podría tener una configuración de dos pares en que más de dos ejecuten un servicio escalable o base de datos paralela aunque ningún nodo esté conectado directamente a los discos que almacenan los datos de la aplicación.

La figura siguiente ilustra una configuración de par en clúster.

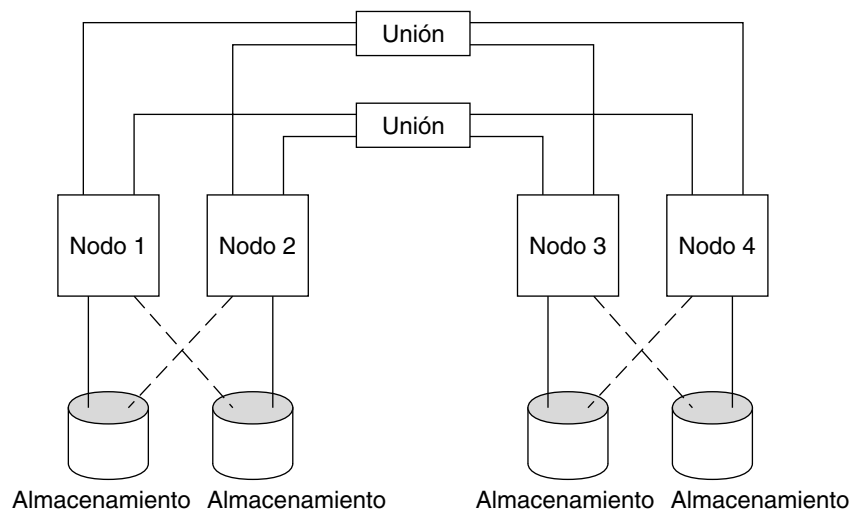


FIGURA 2-2 SPARC: Topología de pares en clúster

SPARC: Topología par+n

La topología par+n incluye un par de nodos conectados directamente al almacenamiento compartido y un conjunto adicional de nodos que usa la interconexión del clúster para acceder al almacenamiento compartido; éstos no tienen conexión directa.

La figura siguiente muestra una topología par+n en que dos de los cuatro nodos (nodo 3 y 4) usan la interconexión del clúster para acceder al almacenamiento. Esta configuración puede ampliarse para que incluya nodos adicionales que no tengan acceso directo al almacenamiento compartido.

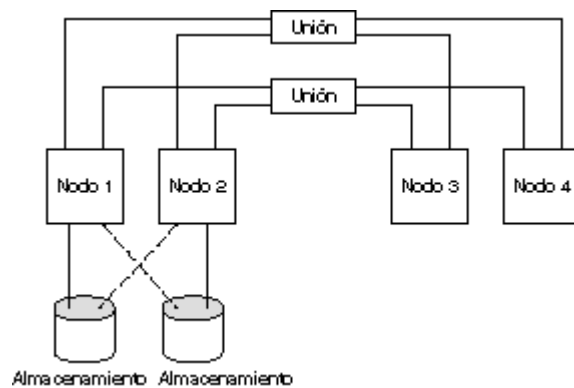


FIGURA 2-3 SPARC: Topología de par+n

SPARC: Topología n+1 (estrella)

Las topologías n+1 incluyen varios nodos primarios y uno secundario que no se deben configurar de forma idéntica. Los nodos primarios proporcionan de forma activa servicios de aplicación. El nodo secundario no debe estar desocupado mientras se espera que uno principal falle.

El nodo secundario es el único de la configuración que está conectado físicamente a todos las unidades de almacenamientos multisistema .

Si se produce un fallo en un nodo primario, Sun Cluster resuelve el fallo transfiriendo los recursos al secundario, dónde éstos funcionan hasta que se conmutan de nuevo (de forma automática o manual) al nodo primario.

El secundario siempre debe tener suficiente capacidad sobrante de CPU para manejar la carga si uno de los primarios falla.

La figura siguiente ilustra la configuración n+1.

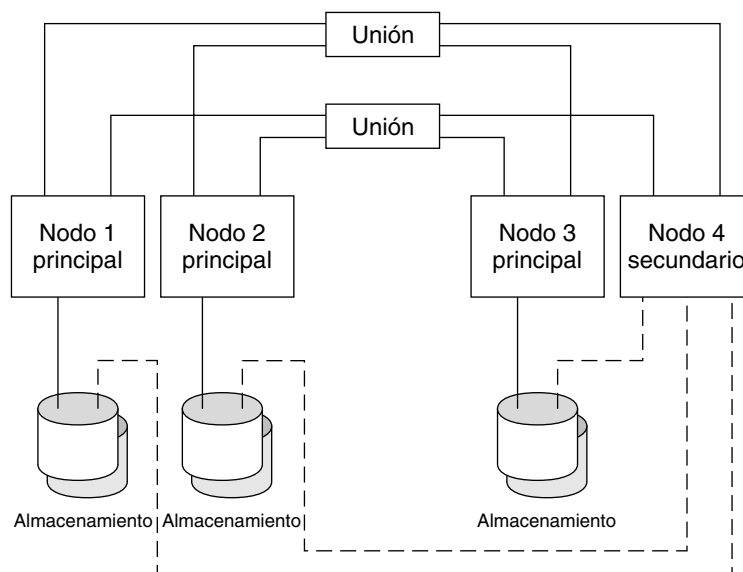


FIGURA 2-4 SPARC: Topología n+1

SPARC: Topología n*n (escalable)

Las topologías n*n permiten a todos los dispositivos de almacenamiento compartido del clúster conectarse con todos los nodos del clúster. Esta topología permite que las aplicaciones de alta disponibilidad se recuperen de fallos de un nodo a otro sin que se produzca una degradación en el servicio. Cuando se produce una recuperación de fallos, el nodo nuevo puede acceder al dispositivo de almacenamiento usando una ruta local en lugar de la interconexión privada.

La siguiente figura muestra una configuración n*n.

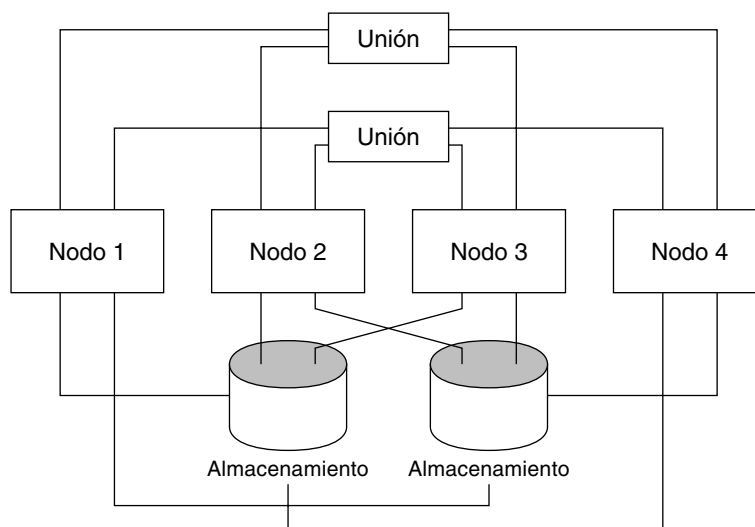


FIGURA 2-5 SPARC: Topología n*n

x86: Ejemplos de topología de Sun Cluster

Una topología es el esquema de conexión que une los nodos del clúster con las plataformas de almacenamiento que se usan en éste. Sun Cluster admite cualquier topología que cumpla las siguientes directrices.

- Sun Cluster, cuando se compone de sistemas basados en plataformas x86, admite dos nodos por clúster.
- Los dispositivos de almacenamiento compartido se deben conectar con ambos nodos.

Sun Cluster no obliga a configurar el clúster mediante topologías específicas. La siguiente topología de par en clúster, la única que admiten los clústers compuestos de nodos basados en plataformas x86, se describe con el único propósito de proporcionar una terminología que explique el esquema de conexiones del clúster. Esta topología es un esquema de conexión típico.

El apartado siguiente incluye una diagrama de ejemplo de la topología.

x86: Topología de par en clúster

Una topología de par en clúster son dos o más pares de nodos que funcionan bajo un único marco administrativo de clúster. En esta configuración sólo se produce recuperación de fallos entre pares. Sin embargo, todos los nodos están conectados por la interconexión del clúster y funcionan bajo el control del software Sun Cluster. Esta topología se podría usar para ejecutar una aplicación de base de datos paralela en un par y una aplicación de recuperación de fallos o de escalabilidad en otro par.

La figura siguiente ilustra una configuración de par en clúster.

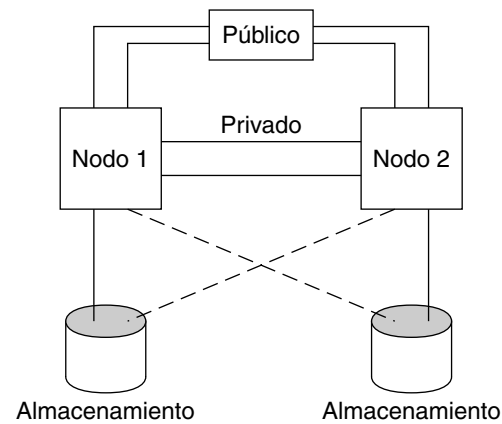


FIGURA 2-6 x86: Topología de par en clúster

Conceptos clave de la administración y desarrollo de las aplicaciones

Este capítulo describe los conceptos clave relacionados con los componentes de hardware de una configuración de sistema de SunPlex. Se tratan los siguientes temas:

- “Interfaces administrativas” en la página 37
- “Hora del clúster” en la página 38
- “Estructura de alta disponibilidad (HA)” en la página 39
- “Dispositivos globales” en la página 41
- “Grupos de dispositivos de discos” en la página 43
- “Espacio de nombres global” en la página 46
- “Sistemas de archivos del clúster” en la página 47
- “Aislamiento de fallos” en la página 55
- “Servicios de datos” en la página 65
- “Desarrollo de nuevos servicios de datos” en la página 72
- “Recursos, grupos de recursos y tipos de recursos” en la página 76
- “Adaptadores de red pública y IP Network Multipathing ” en la página 89
- “SPARC: Compatibilidad con la reconfiguración dinámica” en la página 90

Esta información va dirigida principalmente a administradores del sistema y desarrolladores de aplicaciones que utilicen API y SDK de SunPlex. Los administradores del sistema del clúster puede utilizar esta información para instalar, configurar y administrar el software del clúster. Los desarrolladores de aplicaciones pueden usar la información para entender el entorno del clúster en el que estarán trabajando.

Interfaces administrativas

Se puede elegir como instalar, configurar y administrar el sistema SunPlex desde varias interfaces de usuario, así como realizar tareas de administración del sistema a través de la interfaz gráfica de usuario (GUI) de SunPlex Manager o de la interfaz de línea de órdenes documentada. Además de la interfaz de la línea de órdenes hay

utilidades, como `scinstall` y `scsetup`, que simplifican la instalación seleccionada y las tareas de configuración. El sistema SunPlex también tiene un módulo que se ejecuta como parte de Sun Management Center e incluye una GUI para algunas tareas del clúster, y que está disponible sólo para usarlo en clústers basados en plataformas SPARC. Consulte "Administration Tools" en *Sun Cluster System Administration Guide for Solaris OS* para obtener descripciones completas de las interfaces administrativas.

Hora del clúster

La hora en todos los nodos del clúster debe estar sincronizada. Que esta sincronización se realice con una fuente externa no es importante para el funcionamiento del clúster. El sistema SunPlex emplea Network Time Protocol (NTP) para sincronizar los relojes de los distintos nodos.

En general, una diferencia en el reloj del sistema de una fracción de segundo no causa problemas. Sin embargo, si ejecuta `date (1)`, `rdate (1M)` o `xntpdate (1M)` (interactivamente o desde secuencias `cron`) en un clúster activo, puede forzar un cambio de hora mucho mayor que el de una fracción de segundo para sincronizar el reloj del sistema con el origen de la hora, lo que podría causar problemas con la marca de hora de modificación de archivos o confundir al servicio NTP.

Cuando se instala el sistema operativo Solaris en todos los nodos del clúster, se tiene la oportunidad de cambiar la fecha y hora predeterminadas para el nodo. En general, puede aceptar el valor predeterminado sugerido.

Cuando se instala el software de Sun Cluster mediante `scinstall (1M)`, uno de los pasos del proceso es configurar NTP para el clúster. El software de Sun Cluster incluye un archivo de plantilla, `ntp.cluster` (consulte `/etc/inet/ntp.cluster` en un nodo del clúster instalado), que establece una relación entre iguales entre todos los nodos del clúster, de los cuales uno de ellos es el "preferido". Los nodos se identifican por su nombre de sistema privado y la sincronización de la hora se produce entre la interconexión del clúster. Las instrucciones sobre cómo configurar el clúster para NTP, consulte "Installing and Configuring Sun Cluster Software" en *Sun Cluster Software Installation Guide for Solaris OS*.

También puede configurar uno o más servidores NTP externos al clúster y cambiar el archivo `ntp.conf` para reflejar esta configuración.

Durante el funcionamiento normal, nunca debería ser necesario ajustar la hora en el clúster. Sin embargo, si la hora se ha ajustado incorrectamente cuando instala el sistema operativo Solaris y desea cambiarla, el procedimiento para realizar esta operación se incluye en "Administering the Cluster" en *Sun Cluster System Administration Guide for Solaris OS*.

Estructura de alta disponibilidad (HA)

El sistema SunPlex hace que todos los componentes de la “ruta de acceso” entre usuarios y datos esté muy disponible, incluyendo las interfaces de red, las aplicaciones en sí, el sistema de archivos y los dispositivos multisistema. En general, un componente del clúster está muy disponible si sobrevive a cualquier fallo (de software o hardware) que se produzca en el sistema.

La tabla siguiente muestra los tipos de fallos de componentes de SunPlex (tanto de hardware como de software) y los tipos de recuperación que incorpora la estructura de alta disponibilidad.

TABLA 3-1 Niveles de detección de fallos y recuperación en SunPlex

Componente del clúster fallido	Recuperación de software	Recuperación de hardware
Servicio de datos	API HA , estructura HA	N/D
Adaptador de red pública	IP Network Multipathing	Varias tarjetas adaptadoras de red pública
Sistema de archivos del clúster	Réplicas primaria y secundaria	Dispositivos multisistema
Dispositivo multisistema duplicado	Gestión de volúmenes (Solaris Volume Manager y VERITAS Volume Manager, disponibles sólo para clústers basados en plataformas SPARC)	RAID-5 por hardware (por ejemplo: Sun StorEdge™ A3x00)
Dispositivo global	Réplicas primaria y secundaria	Varias rutas de acceso al dispositivo, uniones de transporte al clúster
Red privada	Software de transporte HA	Varias redes privadas independientes del hardware
Nodo	CMM, controlador de recuperación rápida	Varios nodos

La estructura de alta disponibilidad del software de Sun Cluster detecta el fallo en un nodo rápidamente y crea un servidor equivalente nuevo para los recursos de la estructura en uno de los nodos restantes del clúster. En ningún momento se interrumpe completamente la disponibilidad de los recursos de la estructura. Los que no se hayan visto afectados por el nodo que haya fallado están completamente disponibles durante la recuperación. Además, los recursos de la estructura del nodo fallido vuelven a estar disponibles en cuanto se recuperan. Un recurso de la estructura recuperado no tiene que esperar a que todos los demás recursos de la estructura completen su recuperación.

La mayoría de recursos de la estructura de alta disponibilidad se recuperan de manera transparente para las aplicaciones (servicios de datos) usando el recurso. Las semánticas del acceso a los recursos de la estructura se conservan perfectamente entre los fallos de los nodos. Las aplicaciones simplemente no pueden advertir que el servidor del recurso de la estructura se ha trasladado a otro nodo. El fallo de un único nodo es completamente transparente para los programas desde el resto de los nodos que usan archivos, dispositivos y volúmenes de disco conectados a este nodo, siempre que exista una ruta de acceso alternativa de hardware a los discos desde otro nodo. Un ejemplo es el uso de dispositivos multisistema que tengan puertos con varios nodos.

Supervisor de pertenencia al clúster

Para asegurarse de que los datos permanezcan incorruptos, todos los nodos deben alcanzar un acuerdo uniforme sobre la pertenencia al clúster. Cuando es necesario, CMM coordina una reconfiguración de los servicios del clúster (aplicaciones) en respuesta a un fallo.

CMM recibe información sobre conectividad con otros nodos desde la capa de transporte del clúster. CMM usa la interconexión del clúster para intercambiar información de estado durante la reconfiguración.

Después de detectar un cambio en la composición del clúster, CMM lleva a cabo una configuración sincronizada del clúster en que los recursos de éste podrían redistribuirse de acuerdo con la nueva composición.

A diferencia de anteriores versiones del software Sun Cluster, CMM se ejecuta completamente en el núcleo.

Consulte [“Aislamiento de fallos” en la página 55](#) para obtener información sobre cómo el clúster se protege de una partición en varios clústers independientes.

Mecanismo de recuperación rápida

Si el CMM detecta un problema grave en un nodo, envía una señal a la estructura del clúster para forzar un apagado (aviso grave) del nodo y borrarlo de la pertenencia al clúster. El mecanismo por el que ello ocurre se denomina *recuperación rápida*. Éste obliga a un nodo a apagarse de dos formas.

- Si un nodo deja el clúster e intenta iniciar un clúster nuevo sin tener quórum, se le “aisla” y se le niega el acceso a los discos compartidos. Consulte [“Aislamiento de fallos” en la página 55](#) para obtener detalles sobre este uso de la recuperación rápida.
- Si uno o más daemons específicos del clúster dejan de existir (`clexecd`, `rpc.pmfd`, `rgmd` o `rpc.ed`) CMM detecta el fallo y el nodo entra en pánico. Cuando la finalización de un daemon del clúster hace que un nodo emita un aviso grave, en la consola de éste aparecerá un mensaje parecido a éste.


```
panic[cpu0]/thread=40e60: Failfast: Aborting because "pmfd" died 35 seconds ago.  
409b8 cl_runtime: __0FZsc_syslog_msg_log_no_argsPviTCPcTB+48 (70f900, 30, 70df54, 407acc, 0)  
%10-7: 1006c80 000000a 000000a 10093bc 406d3c80 7110340 0000000 4001 fbf0
```

Después de la condición de aviso grave, el nodo podría reentrar e intentar volver a unirse al clúster o, si éste se compone de sistemas basados en plataformas SPARC, permanecer en el indicador PROM (OBP) de OpenBoot™. La acción que se toma depende del valor del parámetro `auto-boot?`. Puede establecer `auto-boot?` con `eeprom(1M)`, en el indicador `ok` de la PROM de OpenBoot.

Depósito de configuración del clúster (CCR)

CCR usa un algoritmo de puesta al día de dos fases para las actualizaciones: una actualización se tiene que completar satisfactoriamente en todos los miembros del clúster o de lo contrario se anula. CCR usa la interconexión del clúster para aplicar las actualizaciones distribuidas.



Caution – CCR se compone de archivos de texto que nunca se deben editar manualmente, ya que cada archivo contiene un registro de suma de control para asegurar la coherencia entre los nodos. La actualización manual de los archivos de CCR provocaría que un nodo o todo el clúster dejaran de funcionar.

CCR confía en CMM para garantizar que el clúster sólo se ejecute cuando se tenga el suficiente quórum y es responsable de verificar la uniformidad de los datos entre el clúster, efectuando recuperaciones según sea necesario y facilitando actualizaciones a los datos.

Dispositivos globales

El sistema SunPlex usa *dispositivos globales* para ofrecer a todo el clúster un acceso de alta disponibilidad a cualquier dispositivo del clúster, desde cualquier nodo sin que importe dónde esté la conexión física del dispositivo. En general, si un nodo falla mientras se ofrece acceso a un dispositivo global, el software de Sun Cluster descubre automáticamente otra ruta de acceso al dispositivo y redirige el acceso a la misma. Los dispositivos globales de SunPlex pueden ser discos, CD-ROM y cintas. Sin embargo, los discos son los únicos dispositivos globales multipuerto que se admiten. Ello significa que los CD-ROM y los dispositivos de cinta actualmente no son dispositivos de alta disponibilidad. Los discos locales de cada servidor tampoco son multipuerto, por lo tanto no son dispositivos de alta disponibilidad.

El clúster asigna automáticamente ID exclusivas a cada disco, CD-ROM y unidad de cinta del clúster, lo que permite el acceso uniforme a todos los dispositivos desde cualquier nodo del clúster. El espacio de nombres de dispositivos global se conserva en el directorio `/dev/global`. Para obtener más información, véase “Espacio de nombres global” en la página 46.

Los dispositivos globales multipuerto ofrecen más de una ruta de acceso al dispositivo. En el caso de discos multisistema, debido a que forman parte de un grupo de dispositivos de disco alojado por más de un nodo, eso los convierte en discos de alta disponibilidad.

ID de dispositivo (DID)

El software de Sun Cluster gestiona los dispositivos globales a través de una construcción conocida como el pseudocontrolador de ID de dispositivo (DID). Este controlador se utiliza para asignar automáticamente identificadores exclusivos a todos los dispositivos del clúster, incluidos discos multisistema, unidades de cinta y CD-ROM.

El pseudocontrolador de ID de dispositivo (DID) forma parte integral de la prestación de acceso global a los dispositivos del clúster. El controlador de DID examina todos los nodos del clúster y crea una lista de dispositivos de disco única, asignándole a cada uno un número mayor y menor exclusivos que es idéntico en todos los nodos del clúster. El acceso a los dispositivos globales se realiza usando el ID de dispositivo único asignado por el controlador DID en lugar de los ID de dispositivos de Solaris tradicionales, como `c0t0d0` para un disco.

Este enfoque fuerza a que todas las aplicaciones que acceden a discos (como un gestor de volúmenes o aplicaciones que usen dispositivos a bajo nivel) utilicen una ruta de acceso uniforme en todo el clúster. Esta uniformidad es especialmente importante en discos multisistema, debido a que los números mayor y menor de cada dispositivo pueden variar entre distintos nodos, cambiando así también las convenciones de asignación de nombres de dispositivos de Solaris. Por ejemplo, el nodo 1 podría ver un disco multisistema como `c1t2d0` y el nodo2 podría ver ese mismo disco de forma completamente distinta, como `c3t2d0`. El controlador DID asigna un nombre global, como `d10`, que los nodos usarán en su lugar dando a cada uno una correlación uniforme con el disco multisistema.

Los ID de dispositivo se administran con las órdenes `sccidadm(1M)` y `scgdevs(1M)`. Consulte las siguientes páginas de comando man para obtener más información:

- `sccidadm(1M)`
- `scgdevs(1M)`

Grupos de dispositivos de discos

En el sistema SunPlex, todos los dispositivos multisistema deben estar bajo el control del software de Sun Cluster. Primero, cree los grupos de discos de Volume Manager—, mediante conjuntos de discos de Solaris Volume Manager o grupos de discos de VERITAS Volume Manager (disponibles para utilizarlos únicamente en clústeres basados en SPARC)—en los discos multistema. A continuación, se registran los grupos de discos del gestor de discos como *grupos de dispositivos de discos* que forman un tipo de dispositivo global. Además, el software Sun Cluster crea automáticamente un grupo de dispositivos de bajo nivel para cada dispositivo de disco y de cinta del clúster. Sin embargo, estos grupos de dispositivos del clúster permanecen en estado fuera de línea hasta que se accede a ellos como dispositivos globales.

El registro proporciona a SunPlex información del sistema sobre los nodos y sus rutas de acceso a grupos de disco del gestor de volúmenes. En este punto, los grupos de discos del gestor de volúmenes se convierten en accesibles globalmente dentro del clúster. Si hay más de un nodo que pueda escribir (controlar) un grupo de dispositivos de disco, los datos almacenados en éste se consideran de alta disponibilidad. Estos grupos pueden usarse para alojar sistemas de archivos del clúster.

Nota – Los grupos de dispositivos de disco son independientes de los grupos de recursos. Un nodo puede controlar un grupo de recursos (que represente un grupo de procesos de servicio de datos) mientras otro puede controlar los grupos de discos a los que están accediendo los servicios de datos. Sin embargo, la mejor práctica es mantener en el mismo nodo el grupo de dispositivos de disco que almacena los datos de una aplicación determinada y el grupo que contiene sus recursos (el daemon de la aplicación). Consulte “Relationship Between Resource Groups and Disk Device Groups” en *Sun Cluster Data Services Planning and Administration Guide for Solaris OS* para obtener más información acerca de la asociación entre los grupos de recursos y los grupos de dispositivos de disco.

Con un grupo de dispositivos de disco, el grupo de disco del gestor de volúmenes se convierte en “global” porque proporciona soporte de ruta múltiple a los discos que incluye. Cada nodo del clúster conectado físicamente a los discos multisistema proporciona una ruta de acceso al grupo de dispositivos de disco.

Recuperación de fallos del grupo de dispositivos de disco

Debido a que un alojamiento de disco está conectado a más de un nodo, todos los grupos de dispositivos de disco de ese alojamiento estarán disponibles a través de una ruta de acceso alternativa si falla el nodo que esté controlando en ese momento el

grupo de dispositivos. El fallo del nodo que controla el grupo de dispositivos no afecta al acceso al grupo de dispositivos excepto por el tiempo que se tarda en realizar la recuperación y las comprobaciones de integridad. Durante ese tiempo, todas las peticiones se bloquean (de forma transparente para la aplicación) hasta que el sistema vuelve a hacer que el grupo de dispositivos esté disponible.

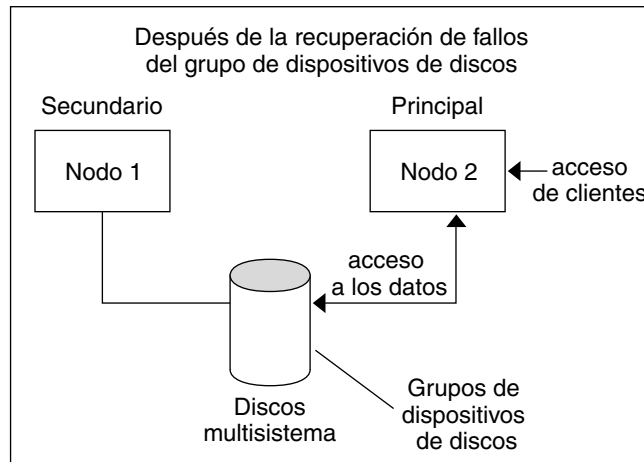
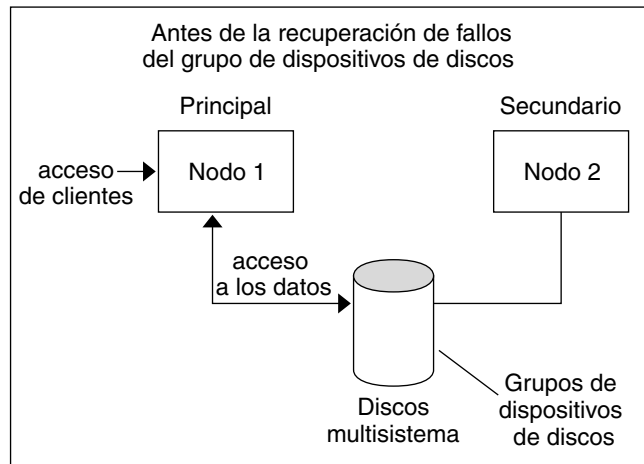


FIGURA 3-1 Recuperación de fallos en grupos de dispositivos de disco

Grupos de dispositivos de disco multipuerto

Este apartado describe las propiedades del grupo de dispositivo de disco que permiten equilibrar el rendimiento y la disponibilidad en una configuración de disco multipuerto. El software Cluster proporciona dos propiedades que se usan para

ajustar una configuración de disco multipuerto: `preferenced` y `numsecondaries`. Con la propiedad `preferenced` se puede controlar el orden en el que los nodos intentan asumir el control cuando se produce una recuperación de fallo. La propiedad `numsecondaries` se utiliza para establecer un número deseado de nodos secundarios para un grupo de dispositivos.

Un servicio de alta disponibilidad se considera caído cuando el nodo primario cae y no hay otros secundarios que puedan promocionar a los primarios. Si se produce la recuperación de fallos y la propiedad `preferenced` es `true`, los nodos siguen el orden de la lista para que se seleccione un secundario. La lista de nodos configurada define el orden en que éstos intentarán asumir el control primario o transicionar de redundante a secundario. Puede cambiar dinámicamente la preferencia de un servicio de dispositivo mediante la utilidad `scsetup(1M)`. La preferencia que está asociada con proveedores de servicio dependientes, por ejemplo un sistema de archivos global, será la del servicio del dispositivo.

El nodo primario comprueba los secundarios durante el funcionamiento normal. En una configuración de disco multipuerto, comprobar todos los nodos secundarios produce una degradación en el rendimiento del clúster y una sobrecarga en la memoria. El soporte para nodos redundantes se ha implementado para minimizar la degradación en el rendimiento y la sobrecarga de memoria que produce el proceso de comprobación. De forma predeterminada el grupo de dispositivos de disco tendrá uno primario y uno secundario. El resto de nodos de proveedor disponibles estarán en línea en el estado redundante. Si se produce una recuperación de fallos, el secundario se convertirá en el primario y el nodo con más prioridad de la lista se convertirá en el secundario.

El número de nodos secundarios que se desee se puede establecer en cualquier número entero entre uno y el número de nodos proveedores no primarios operativos del grupo de dispositivos.

Nota – Si se está utilizando Solaris Volume Manager, se debe crear el grupo de dispositivos de disco antes de establecer la propiedad `numsecondaries` a un número distinto del predeterminado.

De manera predeterminada el número deseado de secundarios para servicios de dispositivos es de uno. El número real de proveedores secundarios que mantiene la estructura de réplica es la deseada, a menos que el número de proveedores no primarios en funcionamiento sea inferior al deseado. Si está añadiendo o quitando nodos de la configuración, deberá hacer cambios en la propiedad `numsecondaries` y revisar bien la lista de nodos. Si se mantiene ésta y el número deseado de secundarios se evitarán conflictos entre el número de secundarios configurado y el número real permitido por la estructura. Utilice la orden `metaset(1M)` para los grupos de dispositivos de Solaris Volume Manager o, si utiliza Veritas Volume Manager, la orden `scconf(1M)` para los grupos de dispositivos de discos VxVM junto con las propiedades `preferenced` y `numsecondaries` para gestionar la adición y supresión

de los nodos de la configuración. Consulte “Administering Cluster File Systems Overview” en *Sun Cluster System Administration Guide for Solaris OS* con el fin de obtener información sobre los procedimientos para cambiar las propiedades de los grupos de dispositivos de discos.

Espacio de nombres global

El mecanismo de software de Sun Cluster que habilita los dispositivos globales es el *espacio de nombres global* que incluye la jerarquía `/dev/global/` así como los espacios de nombres del gestor de volúmenes. El espacio de nombres global representa los discos multisistema y locales (y cualquier otro dispositivo del clúster, como CD-ROM y cintas) e incluye varias rutas de acceso de recuperación de fallos a los discos multisistema. Todos los nodos conectados físicamente a discos multisistema incluyen una ruta de acceso al almacenamiento válida para todos los nodos del clúster.

En general, los espacios de nombres del gestor de volúmenes se encuentran en los directorios `/dev/md/conjunto_discos/dsk` (y `rdsk`), para Solaris Volume Manager, y en los directorios `/dev/vx/dsk/grupo_discos` y `/dev/vx/rdsk/grupo_discos`, para Veritas VxVM. Estos espacios de nombres se componen de directorios para cada conjunto de discos del Solaris Volume Manager y cada grupo de discos de VxVM importados a través del clúster, respectivamente, cada uno de los cuales aloja un nodo de dispositivo para cada metadispositivo o volumen de ese conjunto o grupo de discos.

En el sistema SunPlex todos los nodos de dispositivos del espacio de nombres del gestor de volúmenes local se sustituye por un enlace simbólico con un nodo de dispositivo del sistema de archivos `/global/.devices/node@IDnodo`, donde *IDnodo* es un número entero que representa los nodos del clúster. El software de Sun Cluster sigue presentando los dispositivos del gestor de volúmenes, como enlaces simbólicos, también en sus ubicaciones estándar. Tanto el espacio de nombres global como el estándar del gestor de volúmenes están disponibles desde cualquier nodo del clúster.

Entre las ventajas del espacio de nombres global, se cuentan:

- Cada nodo permanece bastante independiente, con pocos cambios en el modelo de administración de dispositivos.
- Los dispositivos puede convertirse en globales de forma selectiva.
- Los generadores de enlaces de terceros siguen funcionando.
- Dado un nombre de dispositivo local, se ofrece una correlación simple para obtener un nombre global.

Ejemplo de espacios de nombres locales y globales

La tabla siguiente muestra las correlaciones entre los espacios de nombres local y global para un disco multisistema, `c0t0d0s0`.

TABLA 3-2 Correlaciones de espacios de nombres local y global

Componente/Ruta de acceso	Espacio de nombres de nodo local	Espacio de nombres global
Nombre lógico de Solaris	<code>/dev/dsk/c0t0d0s0</code>	<code>/global/.devices/node@ID_nodo/dev/dsk/c0t0d0s0</code>
Nombre DID	<code>/dev/did/dsk/d0s0</code>	<code>/global/.devices/node@ID_nodo/dev/did/dsk/d0s0</code>
Solaris Volume Manager	<code>/dev/md/conjunto-discos/dsk/d0</code>	<code>/global/.devices/node@ID_nodo/dev/md/conjunto_discos/dsk/d0</code>
SPARC: VERITAS Volume Manager	<code>/dev/vx/dsk/grupo_discos/v0</code>	<code>/global/.devices/node@ID_nodo/dev/vx/dsk/grupo_discos/v0</code>

El espacio de nombres global se genera automáticamente durante la instalación y se actualiza con cada arranque de reconfiguración; también se puede generar mediante la orden `scgdevs (1M)`.

Sistemas de archivos del clúster

El sistema de archivos del clúster dispone de las prestaciones siguientes:

- Las ubicaciones de los accesos de archivo son transparentes. Un proceso puede abrir un archivo situado en cualquier parte del sistema y los procesos de todos los nodos pueden usar el mismo nombre de ruta para situar un archivo.

Nota – Cuando el sistema de archivos del clúster lee archivos, no actualiza la hora de acceso en esos archivos.

- Se utilizan protocolos de coherencia para preservar la semántica de acceso a archivos UNIX aunque varios nodos estén accediendo al archivo al mismo tiempo.
- Para mover datos de archivos eficientemente se utiliza masivamente la antememoria y el movimiento de E/S en bloque sin copia.
- El sistema de archivos del clúster ofrece la funcionalidad de bloqueo de archivos a través de las interfaces `fcntl(2)`. Las aplicaciones que se ejecutan en varios nodos del clúster pueden sincronizar el acceso a los datos mediante el bloqueo de

archivo condicional en un sistema de archivos del clúster. Los bloqueos de archivo se recuperan inmediatamente desde los nodos que abandonan el clúster y las aplicaciones que fallan mientras se mantienen los bloqueos.

- El acceso continuo a los datos queda asegurado aunque se produzcan fallos. Las aplicaciones no se ven afectadas por fallos mientras siga estando operativa una ruta de acceso a los discos. Esta garantía se mantiene para el acceso a discos de bajo nivel y todas las operaciones del sistema de archivos.
- Los sistemas de archivos del clúster son independientes del sistema de archivos subyacente y del software de gestión de volúmenes; convierten en global cualquier sistema de archivos en disco admitido.

Se puede montar un sistema de archivos en un dispositivo global con `mount -g` (globalmente) o con `mount` (localmente).

Los programas pueden acceder a los archivos del sistema de archivos del clúster desde cualquier nodo de éste empleando el mismo nombre de archivo (por ejemplo, `/global/foo`).

Los sistemas de archivos del clúster se montan en todos los miembros del clúster, pero no puede montarse en un subconjunto de miembros del clúster.

Los sistemas de archivo clúster no son de tipo diferenciado. Es decir, los clientes ven el sistema de archivos subyacente (por ejemplo, UFS).

Uso de los sistemas de archivos del clúster

En el sistema SunPlex todos los discos multisistema se sitúan en grupos de dispositivos de disco que pueden ser conjuntos de discos del Solaris Volume Manager, grupos de discos de VxVM o discos individuales que no están bajo el control de ningún gestor de volúmenes basado en software.

Para que un sistema de archivos del clúster sea de alta disponibilidad, el almacenamiento en disco subyacente debe estar conectado a más de un nodo. Por consiguiente, un sistema de archivos local (aquel que está almacenado en el disco local de un nodo) que se convierte en un sistema de archivos del clúster no es de alta disponibilidad.

Al igual que ocurre con los sistemas de archivo locales, los sistemas de archivo clúster se pueden montar de dos formas:

- **Manualmente:** con el comando `mount` y las opciones de montaje `-g u` o `-o global` para montar el sistema de archivos del clúster desde la línea de comandos, por ejemplo:

```
SPARC: # mount -g /dev/global/dsk/d0s0 /global/oracle/data
```

- **Automáticamente:** creando una entrada en el archivo `/etc/vfstab` con una opción de montaje `global` para montar el sistema de archivos del clúster desde el arranque. Después puede crear un punto de montaje en el directorio `/global` de

todos los nodos. Éste es una ubicación recomendada, no un requisito. La siguiente es una línea de ejemplo para un sistema de archivos del clúster del archivo `/etc/vfstab`:

```
SPARC: /dev/md/oracle/dsk/d1 /dev/md/oracle/rdisk/d1 /global/oracle/data ufs 2 yes global,logging
```

Nota – Aunque el software de Sun Cluster no impone ninguna política de asignación de nombres a los sistemas de archivos del clúster, puede facilitarse la administración creando un punto de montaje para todos los sistemas de archivos del clúster en el mismo directorio, como `/global/grupo-dispositivo-disco`. Consulte *Sun Cluster 3.1 9/04 Software Collection for Solaris OS (SPARC Platform Edition)* y *Sun Cluster System Administration Guide for Solaris OS* para obtener más información.

Tipo de recurso HASToragePlus

El tipo de recurso HASToragePlus está diseñado para convertir en altamente disponibles configuraciones de sistemas de archivos no globales como UFS y VxFS, permite integrar el sistema de archivos local en el entorno Sun Cluster, convertir aquél en altamente disponible y proporcionar prestaciones del sistema de archivos adicionales, como comprobaciones, montajes y desmontajes forzados que permiten a Sun Cluster recuperarse pasando a sistemas de archivos locales. Para la recuperación de fallos el sistema de archivos local debe residir en grupos de discos globales que tengan habilitados conmutadores de afinidad.

Consulte “Enabling Highly Available Local File Systems” en *Sun Cluster Data Services Planning and Administration Guide for Solaris OS* para obtener información acerca de la forma de utilizar el recurso HASToragePlus.

Éste también se puede usar para sincronizar el inicio de recursos y grupos de dispositivos de disco de los que dependen los recursos. Para obtener más información, consulte “Recursos, grupos de recursos y tipos de recursos” en la página 76.

Opción de montaje syncdir

La opción de montaje `syncdir` puede utilizarse en sistemas de archivos del clúster que usen UFS como sistema de archivo subyacente. Sin embargo, existe una mejora de rendimiento significativa si no se especifica `syncdir`. Si se especifica `syncdir`, se garantiza que las escrituras sean compatibles con POSIX. Si no se especifica, tendrá el mismo comportamiento que se ve con sistemas de archivos NFS. Por ejemplo, en algunos casos sin `syncdir`, no se descubriría una condición de falta de espacio hasta que se cerrara el archivo. Con `syncdir` (y el comportamiento POSIX), la condición de falta de espacio se descubriría durante la operación de escritura. La posibilidad de tener problemas si no se especifica `syncdir` es remota, por ello es recomendable no especificarla ya que así se mejora el rendimiento.

Si usa un clúster basado en la plataforma SPARC, Veritas VxFS no tiene una opción de montaje equivalente a `syncdir` para UFS. El comportamiento de VxFS es el mismo que para UFS cuando no se especifica la opción de montaje `syncdir`.

Consulte [“FAQ sobre sistemas de archivos” en la página 96](#) para obtener preguntas frecuentes sobre los dispositivos globales y los sistemas de archivos del clúster.

Supervisión de las rutas de disco

La versión actual del software Sun Cluster admite la supervisión de las rutas del disco (DPM). Este apartado incluye información conceptual sobre DPM, el daemon DPM, y las herramientas de administración que se pueden usar para supervisar las rutas del disco. Consulte *Sun Cluster 3.1 9/03 System Administration Guide* para obtener información de los procedimientos para supervisar, dejar de supervisar y comprobar el estado de las rutas del disco.

Nota – DPM no se admite en nodos que ejecuten versiones anteriores al Software de Sun Cluster 3.1 4/04. No utilice órdenes de DPM durante una modernización. Una vez finalizada la modernización en todos los nodos, éstos deben estar en línea para poder utilizar las órdenes de DPM.

Información general

DPM mejora la disponibilidad general de la recuperación de fallos y la conmutación por supervisión de la disponibilidad de la ruta de acceso a disco secundaria. Utilice el comando `scdpm` para verificar la disponibilidad de la ruta del disco que utiliza un recurso antes de conmutarlo. Las opciones que se incluyen con la orden `scdpm` permiten supervisar las rutas de acceso a discos hacia un nodo individual o hacia todos los nodos del clúster. Consulte la página de comando `man scdpm(1M)` para obtener más información sobre las opciones de la línea de órdenes.

Los componentes DPM se instalan a partir del paquete `SUNWscu`. Éste lo instala el procedimiento de instalación estándar de Sun Cluster. Consulte la página de comando `man scinstall(1M)` para obtener detalles sobre la instalación de la interfaz. La tabla siguiente describe la ubicación predeterminada de los componentes DPM.

Ubicación	Componente
Daemon	<code>/usr/cluster/lib/sc/scdpmd</code>

Ubicación	Componente
Interfaz de línea de órdenes	/usr/cluster/bin/scdpm
Bibliotecas compartidas	/user/cluster/lib/libscdpm.so
Archivo de estado de daemon (creado en tiempo de ejecución)	/var/run/cluster/scdpm.status

En cada nodo se ejecuta un daemon DPM de subprocesso múltiple. El daemon DPM (`scdpmd`) lo inicia la secuencia `rc.d` cuando arrancan los nodos. Si surge algún problema, `pmfd` gestiona el daemon y lo reinicia automáticamente. La lista siguiente describe cómo funciona `scdpmd` en el arranque inicial.

Nota – En el arranque, el estado de cada ruta del disco se inicializa a UNKNOWN.

1. El daemon DPM recoge información de rutas del disco y nombres de nodo del archivo de estado anterior o de la base de datos CCR. Para obtener más información sobre CCR, consulte [“Depósito de configuración del clúster \(CCR\)” en la página 41](#). Cuando se inicia el daemon DPM, éste puede forzarse para que lea la lista de discos supervisados a partir de un nombre de archivo especificado.
2. El daemon DPM inicializa la interfaz de comunicaciones para que responda a solicitudes de componentes que son externos al mismo, como la interfaz de línea de órdenes.
3. El daemon DPM realiza un ping a cada ruta del disco de la lista supervisada cada 10 minutos mediante órdenes de tipo `scsi_inquiry`. Todas las entradas están bloqueadas para evitar que la interfaz de comunicaciones acceda al contenido de una entrada que esté siendo modificada.
4. El daemon DPM envía una notificación a Sun Cluster Event Framework y registra el estado nuevo de la ruta a través del mecanismo UNIX `syslogd(1M)`.

Nota – Todos los errores relacionados con el daemon se tratan en `pmfd(1M)`. Todas las funciones de la API devuelven 0 cuando tienen éxito y -1 cuando se produce algún error.

El daemon DPM supervisa la disponibilidad de la ruta lógica que es visible a través de controladores de ruta múltiple como MPxIO, HDLM y PowerPath. Las rutas de acceso físicas individuales gestionadas por estos controladores no están supervisadas porque el controlador de ruta múltiple oculta los fallos individuales al daemon DPM.

Supervisión de las rutas del disco

Este apartado describe dos métodos para supervisar las rutas del disco en el clúster. El primer método lo ofrece el comando `scdpm`. Utilice esta orden para supervisar, dejar de supervisar o mostrar el estado de las rutas del disco del clúster; también resulta útil para imprimir la lista de discos fallidos y supervisar rutas de los discos a partir de un archivo.

El segundo método para supervisar las rutas de los discos en el clúster lo ofrece la interfaz gráfica de usuario (GUI) de SunPlex Manager. SunPlex Manager incluye una vista topográfica de las rutas del disco supervisadas del clúster. La vista se actualiza cada 10 minutos para incluir información sobre el número de pings que han fallado. Para administrar rutas del disco use la información que proporciona la GUI de SunPlex Manager junto con la orden `scdpm(1M)`. Consulte “Administering Sun Cluster With the Graphical User Interfaces” en *Sun Cluster System Administration Guide for Solaris OS* para obtener información acerca de SunPlex Manager.

Uso de la orden `scdpm` para supervisar las rutas del disco

El comando `scdpm(1M)` proporciona a DPM comandos de administración que permiten realizar las tareas siguientes:

- Supervisar una ruta del disco nueva
- Dejar de supervisar una ruta del disco
- Volver a leer los datos de configuración de la base de datos CCR
- Leer los discos para supervisar o dejar de hacerlo a partir de un archivo especificado
- Informar del estado de una o todas las rutas de acceso de disco del clúster
- Imprimir todas las rutas de acceso de disco que son accesibles desde un nodo

Emita la orden `scdpm(1M)` con el argumento `ruta-disco` desde cualquier nodo para llevar a cabo tareas de administración de DPM en el clúster. Éste consta en todos los casos de un nombre de nodo y de un nombre de disco. El primer nodo no es necesario y, en caso de no especificarlo, toma el valor `all`. La tabla siguiente describe las convenciones de asignación de nombres para la ruta del disco.

Nota – Se recomienda utilizar nombres de ruta del disco globales, ya que son coherentes dentro de todo el clúster. Los nombres de las rutas del disco UNIX no son coherentes en todo el clúster. La ruta del disco UNIX correspondiente a un disco determinado puede diferir en distintos nodos del clúster. La ruta podría ser `c1t0d0` en un nodo y `c2t0d0` en otro. Si utiliza nombres de ruta del disco UNIX, utilice la orden `scdidadm -L` para asignar el nombre UNIX al nombre global antes de ejecutar órdenes de DPM. Consulte la página de comando `man scdidadm(1M)`.

TABLA 3-3 Ejemplos de nombres de rutas del disco

Tipo de nombre	Ejemplo de nombre de ruta del disco	Descripción
Ruta del disco global	<code>schost-1:/dev/did/dsk/d1</code>	Ruta del disco d1 en el nodo <code>schost-1</code>
	<code>all:d1</code>	Ruta del disco d1 en todos los nodos del clúster
Ruta del disco UNIX	<code>schost-1:/dev/rdisk/c0t0d0s0</code>	Ruta del disco <code>c0t0d0s0</code> en el nodo <code>schost-1</code>
	<code>schost-1:all</code>	Todas las rutas del nodo <code>schost-1</code>
Todas las rutas del disco	<code>all:all</code>	Todas las rutas del disco de todos los nodos del clúster

Uso de SunPlex Manager para supervisar las rutas del disco

SunPlex Manager permite realizar las siguientes tareas de administración DPM básicas:

- Supervisar una ruta del disco
- Dejar de supervisar una ruta del disco
- Visualizar el estado de todas las rutas del disco del clúster.

Consulte la ayuda en línea de SunPlex Manager para obtener información de procedimientos para realizar una administración de la ruta del disco utilizando SunPlex Manager.

Quórum y dispositivos de quórum

Esta sección se divide en los siguientes apartados:

- “Acerca de los recuentos de votos de quórum” en la página 55
- “Aislamiento de fallos” en la página 55
- “Acerca de las configuraciones de quórum” en la página 57
- “Cumplimiento de los requisitos de dispositivos de quórum ” en la página 58
- “Cumplimiento de las mejores prácticas recomendadas de dispositivos de quórum” en la página 58

- “Configuraciones de quórum recomendadas ” en la página 60
- “Configuraciones de quórum atípicas ” en la página 63
- “Configuraciones de quórum malas ” en la página 64

Nota – Para obtener una lista de los dispositivos específicos que Sun Cluster admite como dispositivos de quórum, póngase en contacto con el proveedor de servicios Sun.

Debido a que los nodos del clúster comparten datos y recursos, un clúster nunca se debe dividir en particiones separadas que estén activas a la vez porque varias particiones activas pueden provocar que se dañen los datos. Cluster Membership Monitor (CMM) y el algoritmo de quórum garantizan que como mucho una instancia del mismo clúster está operativa cada vez, incluso si se particiona la interconexión del clúster.

Si desea obtener más información acerca de CMM, consulte “Cluster Membership” en *Sun Cluster Overview for Solaris OS*

Pueden surgir dos tipos de problemas derivados de las particiones del disco:

- Esquizofrenia
- Amnesia

La esquizofrenia se produce cuando se pierde la interconexión del clúster entre nodos y el clúster se particiona en clústeres secundarios. Cada partición cree que es la única partición porque los nodos en una partición no se pueden comunicar con los nodos en otra partición.

La amnesia se produce cuando el clúster se reinicia después de un apagado teniendo datos de configuración del clúster más antiguos que los del momento del apagado. Este problema se puede producir cuando inicia el clúster en un nodo que no se encuentra en la partición que estaba funcionando la última vez.

Sun Cluster evita la esquizofrenia y amnesia:

- Asignando un voto a cada nodo
- Controlando la mayoría de los votos de un clúster operativo

Una partición con la mayoría de votos tiene *quórum* y se le permite funcionar. Este mecanismo de votos de la mayoría evita la esquizofrenia y amnesia cuando hay más de dos nodos configurados en un clúster. Sin embargo, el recuento de los votos de los nodos por sí solo no es suficiente cuando hay más de dos nodos configurados en un clúster. En un clúster de dos nodos, la mayoría es dos. Si dicho clúster de dos nodos se particiona, es necesario un voto externo para que cualquiera de las particiones obtenga quórum. Este voto externo lo proporciona un *dispositivo del quórum*

Acerca de los recuentos de votos de quórum

Use la opción `-q` del comando `scstat` para determinar la siguiente información:

- Votos totales configurados
- Votos presentes actuales
- Votos necesarios para quórum

Para obtener más información acerca de este comando, consulte `scstat(1M)`.

Los nodos y los dispositivos de quórum aportan votos al clúster para formar quórum.

Un nodo aporta votos en función del estado del nodo:

- Un nodo tiene un recuento de votos de *uno* cuando arranca y es miembro del clúster.
- Un nodo tiene un recuento de voto de *cero* cuando el nodo se está instalando.
- Un nodo tiene un recuento de voto de *cero* cuando un administrador de sistema pone el nodo en estado de mantenimiento.

Los dispositivos de quórum aportan votos basándose en el número de votos conectados al dispositivo. Cuando configura un dispositivo de quórum, Sun Cluster asigna al dispositivo de quórum un recuento de votos de $N-1$ donde N es el número de votos conectados al dispositivo de quórum. Por ejemplo, un dispositivo de quórum conectado a dos nodos con recuentos distintos de cero, tiene un recuento del quórum igual a uno (dos menos uno).

Un dispositivo de quórum aporta votos si se cumple *una* de las siguientes condiciones:

- Al menos uno de los nodos a los que está conectado actualmente el dispositivo de quórum es miembro del clúster.
- Al menos uno de los nodos a los que el dispositivo de quórum está conectado está arrancando y dicho nodo era miembro de la última partición de clúster propietaria del dispositivo de quórum.

Puede configurar los dispositivos de quórum durante la instalación del clúster o posteriormente utilizando los procedimientos que se describen en “Administering Quorum” en *Sun Cluster System Administration Guide for Solaris OS*.

Aislamiento de fallos

Un problema fundamental de los clústers es un fallo que provoque en éstos una partición (denominada *esquizofrenia*). Cuando ocurre, no todos los nodos pueden comunicarse, por lo que algunos podrían intentar formar clústers individuales o subconjuntos que se crearían con permisos de acceso y de propiedad exclusivos respecto a los discos multisistema. En consecuencia, varios nodos intentando escribir en los discos podrían provocar la corrupción de los datos.

El aislamiento de fallos limita el acceso de los nodos a los dispositivos multisistema, evitando que físicamente se pueda acceder a ellos. Cuando un nodo abandona el clúster (falla o se particiona), el aislamiento de fallos se asegura de que el nodo ya no pueda acceder a los discos. Sólo los nodos miembros actuales tendrán acceso a los discos, conservándose así la integridad de los datos.

Los servicios de dispositivos de disco ofrecen prestaciones para servicios que utilizan dispositivos multisistema. Cuando un miembro del clúster que sirve actualmente como primario (propietario) del grupo de dispositivos de disco cae o deja de ser accesible, se elige otro primario, lo que permite que continúe el acceso al grupo de dispositivos de disco con la mínima interrupción. Durante este proceso, el primario antiguo debe renunciar al acceso a los dispositivos antes de que el nuevo primario pueda iniciarse. Sin embargo, cuando un miembro se descuelga del clúster y deja de estar disponible, éste no puede informar al nodo que libere los dispositivos para los que era primario. Por tanto, se necesita un medio para permitir que los miembros supervivientes tomen control y accedan a los dispositivos globales de los miembros fallidos.

El sistema SunPlex usa reservas de disco SCSI para implementar el aislamiento de fallos, gracias a las cuales, los nodos fallidos se “aislan” de los dispositivos multisistema, evitando que accedan a estos discos.

Las reservas de los discos SCSI-2 admiten un tipo que concede acceso a todos los nodos conectados al disco (cuando no hay ninguna reserva vigente) o restringe el acceso a un nodo individual (el nodo que retiene la reserva).

Cuando un miembro del clúster detecta que otro nodo ya no se está comunicando a través de la interconexión del clúster, inicia un procedimiento de aislamiento de fallos para evitar que el otro nodo acceda a los discos compartidos. Cuando se produce este aislamiento de fallos es normal que el nodo aislado entre en pánico con mensajes de “conflicto de reserva” en la consola.

El conflicto de reserva se produce porque después de haberse detectado que el nodo ya no es miembro del clúster, se pone una reserva SCSI en todos los discos que están compartidos entre este nodo y los demás nodos. El nodo podría no advertir que se le está aislando y si intenta acceder a uno de los discos compartidos, detecta la reserva y entra en situación de pánico.

Mecanismo de recuperación rápida para aislamiento de fallos

El mecanismo por el que la estructura del clúster se asegura de que un nodo fallido no pueda reanunciar y empezar a escribir en almacenamiento compartido se denomina *recuperación rápida*.

Los nodos que son miembros del clúster habilitan permanentemente un ioctl específico, `MHIOCENFAILFAST`, para los discos a los que tienen acceso, incluidos los de quórum. `ioctl` es una directiva para el controlador de disco y da a un nodo la posibilidad de entrar en pánico por sí mismo si éste no puede acceder al disco debido a que otro nodo lo está reservando.

`ioctl MHIOCENFAILFAST` hace que el controlador compruebe el retorno del error de cada lectura y escritura que emite un nodo al disco en el caso del código de error `Reservation_Conflict`. `ioctl` emite periódicamente en segundo plano una operación de prueba al disco para comprobar si se produce `Reservation_Conflict`. Las rutas del flujo de control de segundo y primer plano entran en pánico si se devuelve `Reservation_Conflict`.

En discos SCSI-2, las reservas no son persistentes, pues no resisten los rearranques de los nodos. En los discos SCSI-3 con reserva de grupo persistente (PGR), la información de reserva se almacena en el disco y permanece entre los rearranques de los nodos. El mecanismo de recuperación rápida funciona igual aunque disponga de discos SCSI-2 o SCSI-3.

Si un nodo pierde conectividad con los otros nodos del clúster y no es parte de una partición que pueda conseguir el quórum, otro nodo lo expulsa del clúster. Otro nodo que forme parte de la partición que pueda conseguir el quórum coloca reservas en los discos compartidos y cuando el nodo que no tiene el quórum intenta acceder a éstos recibe un conflicto de reserva y entra en condición de pánico como resultado del mecanismo de recuperación rápida.

Después de la condición de aviso grave, el nodo podría rearrancar e intentar volver a unirse al clúster o, si éste se compone de sistemas basados en plataformas SPARC, permanecer en el indicador PROM (OBP) de OpenBoot™. La acción que se toma depende del valor del parámetro `auto-boot?`. Se puede establecer `auto-boot?` con `eeeprom(1M)`, en el indicador `ok` de la PROM de OpenBoot en un clúster basado en la plataforma SPARC o con la utilidad SCSI que se puede ejecutar opcionalmente tras un arranque de la BIOS en un clúster basado en la plataforma x86.

Acerca de las configuraciones de quórum

La siguiente lista contiene hechos acerca de las configuraciones de quórum:

- Los dispositivos de quórum pueden contener datos de usuario.
- En una configuración N+1 donde *N* dispositivos de quórum están conectados a uno de los nodos de 1 a *N* y al nodo *N*+1, el clúster sobrevive a la desactivación de los nodos 1 a *N* o cualquiera de los nodos *N*/ 2. Esta disponibilidad asume que el dispositivo de quórum está funcionando correctamente.
- En una configuración de *N* nodos en la que un único dispositivo de quórum se conecta a todos los nodos, el clúster puede sobrevivir a la desconexión de cualquiera de los nodos *N*- 1. Esta disponibilidad asume que el dispositivo de quórum está funcionando correctamente.

- En una configuración de N nodos donde un único dispositivo de quórum se conecta a todos los nodos, el clúster puede sobrevivir al fallo del dispositivo de quórum si todos los nodos del clúster están disponibles.

Para obtener ejemplos de configuraciones de quórum que se deben evitar, consulte “Configuraciones de quórum malas ” en la página 64. Para obtener ejemplos de configuraciones de quórum recomendadas, consulte “Configuraciones de quórum recomendadas ” en la página 60.

Cumplimiento de los requisitos de dispositivos de quórum

Debe cumplir los siguientes requisitos. En caso contrario, puede poner en peligro la disponibilidad del clúster.

- Asegúrese de que Sun Cluster admite el dispositivo específico como dispositivo de quórum.

Nota – Para obtener una lista de los dispositivos específicos que Sun Cluster admite como dispositivos de quórum, póngase en contacto con el proveedor de servicios Sun.

Sun Cluster admite dos tipos de dispositivos de quórum:

- Discos compartidos multisistema que admiten reservas SCSI-3 PGR
- Discos compartidos de doble sistema que admiten reservas SCSI-2
- En una configuración de dos nodos, debe configurar al menos un dispositivo de quórum para asegurar que un único nodo puede continuar si el otro nodo falla. Consulte [Figura 3–2](#).

Para obtener ejemplos de configuraciones de quórum que se deben evitar, consulte “Configuraciones de quórum malas ” en la página 64. Para obtener ejemplos de configuraciones de quórum recomendadas, consulte “Configuraciones de quórum recomendadas ” en la página 60.

Cumplimiento de las mejores prácticas recomendadas de dispositivos de quórum

Use la siguiente información para evaluar la mejor configuración de quórum para su topología:

- ¿Tiene un dispositivo capaz de estar conectado a todos los nodos del clúster?



- En caso afirmativo, configure dicho dispositivo como el dispositivo de quórum. *No es necesario configurar otro dispositivo de quórum porque la configuración es la configuración óptima.*

Precaución – Si ignora este requisito y añade otro dispositivo de quórum, el dispositivo de quórum adicional reduce la disponibilidad del clúster.

- En caso contrario, configure el dispositivo(s) de doble puerto.
- Asegúrese de que el número de votos aportados por los dispositivos de quórum es menor que el número de votos total aportados por los nodos. En caso contrario los nodos no pueden formar un clúster si todos los discos no están disponibles—, incluso si todos los nodos están funcionando.

Nota – En ocasiones, en entornos particulares, puede ser preferible reducir la disponibilidad general del clúster para satisfacer sus necesidades. En estas situaciones, puede ignorar estas recomendaciones. Sin embargo, no cumplir estas recomendaciones reduce la disponibilidad general. Por ejemplo, en la configuración que se define en [“Configuraciones de quórum atípicas” en la página 63](#) el clúster tiene una menor disponibilidad: los votos del quórum superan los votos del nodo. El clúster tiene la propiedad de que si se pierde el acceso al almacenamiento compartido entre los Nodos A y B, todo el clúster fallará.

Consulte [“Configuraciones de quórum atípicas” en la página 63](#) para conocer las excepciones a estas recomendaciones de mejores prácticas.

- Especifique un dispositivo de quórum entre cada par de nodos que compartan el acceso al dispositivo de almacenamiento. Esta configuración de quórum acelera el proceso de aislamiento de fallos. Consulte [“Quórum en configuraciones de más de dos nodos” en la página 61](#).
- En general, si la adición de un dispositivo de quórum iguala el número total de votos del clúster, la disponibilidad del clúster disminuye.
- Los dispositivos de quórum ralentizan ligeramente las reconfiguraciones después de que se una un nodo nuevo o se desactive uno antiguo. Por tanto, no agregue más dispositivos de quórum de los necesarios.

Para obtener ejemplos de configuraciones de quórum que se deben evitar, consulte [“Configuraciones de quórum malas” en la página 64](#). Para obtener ejemplos de configuraciones de quórum recomendadas, consulte [“Configuraciones de quórum recomendadas” en la página 60](#).

Configuraciones de quórum recomendadas

Para obtener ejemplos de configuraciones de quórum que se deben evitar, consulte "Configuraciones de quórum malas" en la página 64.

Quórum en configuraciones de dos nodos

Son necesarios dos votos de quórum para que se forme un clúster de dos nodos, que pueden provenir de los dos nodos del clúster o de un nodo y un dispositivo del quórum.

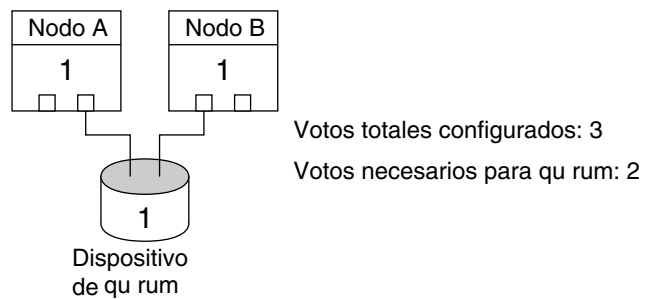
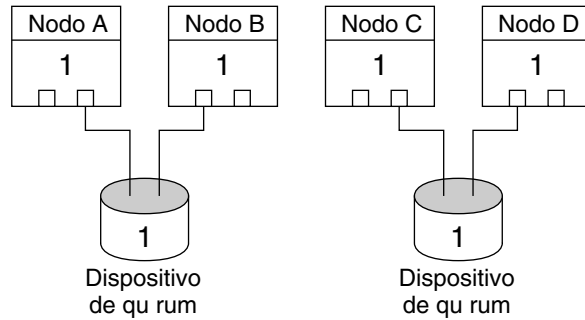


FIGURA 3-2 Configuración de dos nodos

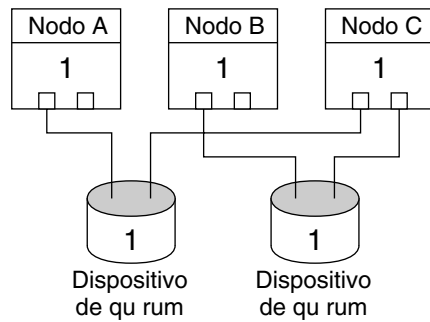
Quórum en configuraciones de más de dos nodos

Es válido configurar un clúster de más de dos nodos sin dispositivo de quórum. Sin embargo, si realiza esta operación, no podrá iniciar el clúster sin una mayoría de nodos en el clúster.



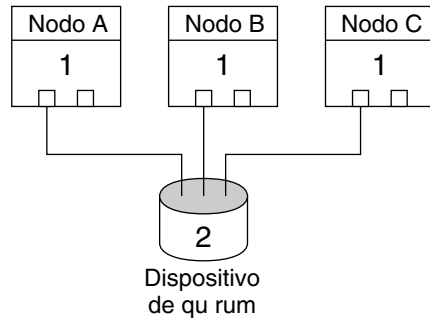
Votos totales configurados: 6
 Votos necesarios para quorum: 4

En esta configuración, cada par debe estar disponible para cada par para sobrevivir.



Votos totales configurados: 5
 Votos necesarios para quorum: 3

En esta configuración, las aplicaciones se configuran normalmente para ejecutarse en el Nodo A y el Nodo B y utilizar el Nodo C como respaldo.



Votos totales configurados: 5
 Votos necesarios para quorum: 3

En esta configuración, la combinación de uno o más nodos y el dispositivo de quorum pueden formar un clúster.

Configuraciones de quórum atípicas

Figura 3-3 asume que se están ejecutando aplicaciones de misión crítica (una base de datos de Oracle por ejemplo) en el Nodo A y Nodo B. Si el Nodo A y el Nodo B no están disponibles y no se puede acceder a los datos compartidos, es posible que prefiera que todo el clúster se apague. En caso contrario, esta configuración no es óptima porque no proporciona una alta disponibilidad.

Para obtener más información sobre las prácticas más recomendadas para esta excepción, consulte “Cumplimiento de las mejores prácticas recomendadas de dispositivos de quórum” en la página 58.

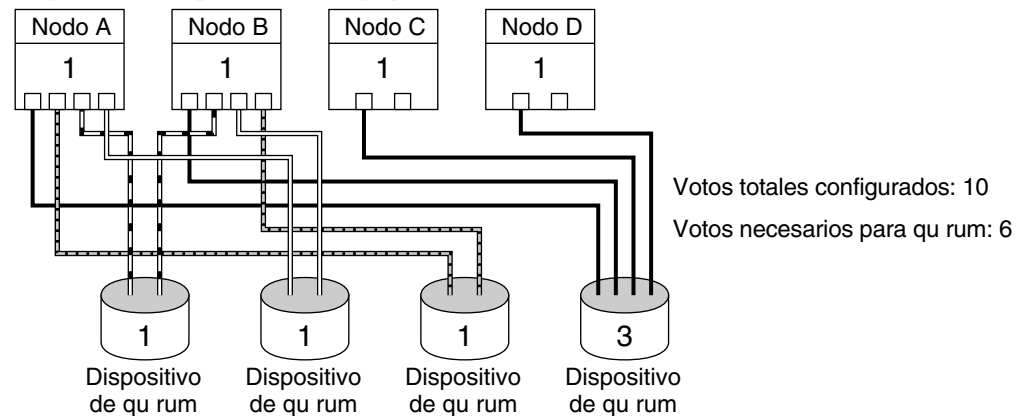
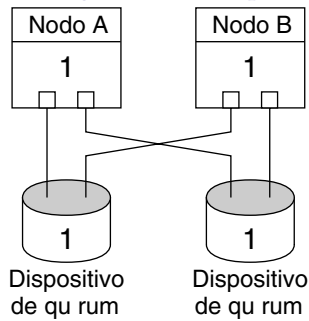


FIGURA 3-3 Configuración atípica

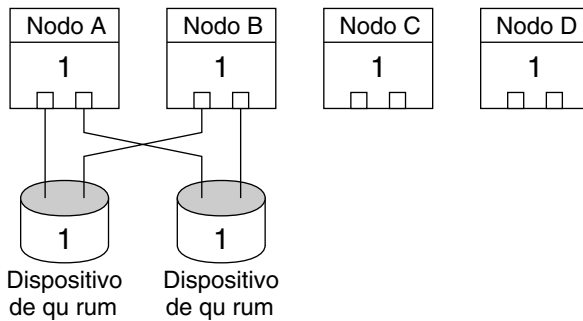
Configuraciones de quórum malas

Para obtener ejemplos de configuraciones de quórum recomendadas, consulte ["Configuraciones de quórum recomendadas"](#) en la página 60.



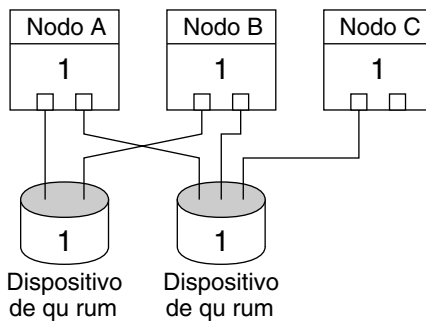
Votos totales configurados: 4
Votos necesarios para quórum: 3

Esta configuración viola las prácticas recomendadas que indican que los votos de dispositivo de quórum deben ser estrictamente menos que los votos de los nodos.



Votos totales configurados: 6
Votos necesarios para quórum: 4

Esta configuración viola las prácticas recomendadas que indican que no se deben agregar dispositivos de quórum para igualar los votos totales. Esta configuración no añade disponibilidad.



Votos totales configurados: 5
Votos necesarios para quórum: 3

Esta configuración viola las prácticas recomendadas que indican que los votos de dispositivo de quórum deben ser estrictamente menos que los votos de los nodos.

Servicios de datos

El término *servicio de datos* describe una aplicación de otras fabricantes, como Sun Java System Web Server (anteriormente Sun Java System Web Server) o, en el caso de clústers basados en plataformas SPARC, Oracle, configurada para ejecutarse en un clúster antes que en un único servidor. Un servicio de datos se compone de una aplicación, archivos de configuración de Sun Cluster y métodos de gestión Sun Cluster que controlan las acciones siguientes de la aplicación.

- Start
- Detener
- Supervisar y tomar acciones correctoras
- Para obtener información acerca de los servicios de datos, consulte “Data Services” en *Sun Cluster Overview for Solaris OS*.

En la [Figura 3-4](#) se compara una aplicación que se ejecuta en un servidor de aplicaciones individual (el modelo de servidor individual) con la misma aplicación que se ejecuta en un clúster (el modelo de servidores en clúster). Tenga en cuenta que desde el punto de vista del usuario, no hay diferencia entre las dos configuraciones, excepto en que la aplicación en clúster podría ejecutarse más rápido y con alta disponibilidad.

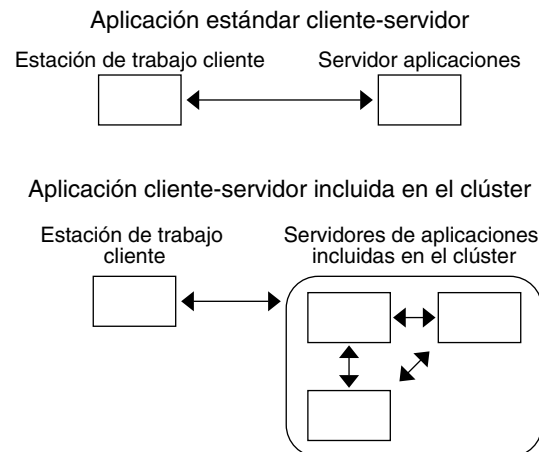


FIGURA 3-4 Configuración cliente/servidor estándar frente a configuración en clúster

En el modelo de servidor individual la aplicación se configura para que acceda al servidor a través de una interfaz de red pública determinada (un nombre de sistema). El nombre de sistema está asociado con este servidor físico.

En el modelo basado en clúster, la interfaz de red pública es un *nombre de sistema lógico* o una *dirección compartida*. El término *recursos de red* se utiliza para referirse a los nombres de sistema lógicos y para las direcciones compartidas.

Algunos servicios de datos obligan a especificar nombres de sistema lógicos o direcciones compartidas como interfaces de red, no son intercambiables; otros, en cambio, permiten especificar nombres de sistema lógicos o direcciones compartidas. Consulte la instalación y configuración de cada servicio de datos para obtener los detalles que se deben especificar.

Un recurso de red no está asociado con ningún servidor físico determinado, puede cambiar entre servidores físicos.

Un recurso de red está asociado inicialmente con un nodo, el *primario*. Si éste falla, los recursos de la red y de la aplicación se trasladan a otro nodo del clúster (uno secundario). Cuando el recurso de red se recupera del problema, después de un breve intervalo, el recurso de aplicación continúa ejecutándose en el secundario.

En la [Figura 3-5](#) se compara el modelo de servidor individual con el basado en clúster. Tenga en cuenta que en el modelo de servidor basado en clúster un recurso de red (nombre de sistema lógico, en este ejemplo) puede trasladarse entre dos o más nodos del clúster. La aplicación está configurada para usar este nombre de sistema lógico en lugar de uno asociado a un servidor en particular.

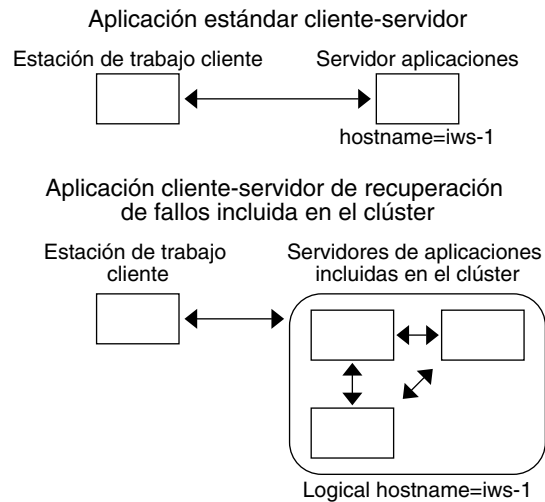


FIGURA 3-5 Comparación entre los nombres de sistema fijo y lógico

Una dirección compartida también está asociada inicialmente con un nodo denominado nodo de interfaz global (GIF). La dirección compartida se usa como interfaz de red única en el clúster que se conocen como *interfaces globales*.

La diferencia entre los modelos de nombre de sistema lógico y de servicio escalable es que en éste cada nodo también tiene la dirección compartida configurada activamente en su interfaz de bucle. Esta configuración hace posible tener activas varias instancias de un servicio de datos en varios nodos simultáneamente. El término “servicio escalable” significa que puede agregarse más potencia de CPU a la aplicación agregando nodos del clúster adicionales con lo que el rendimiento se multiplicará.

Si el nodo GIF falla, la dirección compartida puede llevarse a otro nodo que también esté ejecutando una instancia de la aplicación (convirtiendo así al otro nodo en el nuevo GIF). O, la dirección compartida puede trasladarse a otro nodo del clúster que no estuviera ejecutando la aplicación previamente.

En la [Figura 3-6](#) se compara la configuración de servidor individual con la configuración de servicio escalable en clúster. Tenga en cuenta que, en la configuración de servicio escalable, la dirección compartida está presente en todos los nodos. De forma análoga a como se usan los nombres de sistema para los servicios de datos de recuperación de fallos, la aplicación se configura para usar esta dirección compartida en lugar de un nombre de sistema asociado a un servidor determinado.

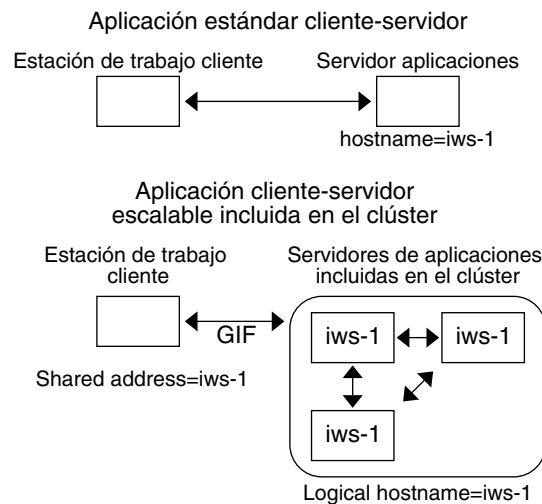


FIGURA 3-6 Comparación entre los nombres de sistema fijos y las direcciones compartidas

Métodos de servicios de datos

El software Sun Cluster proporciona un conjunto de métodos de gestión de servicios que se ejecutan bajo el control del Resource Group Manager (RGM), el cual los usa para iniciar, detener y supervisar la aplicación en los nodos del clúster. Estos métodos, junto con el software de la estructura del clúster y los dispositivos multisistema, permiten a las aplicaciones convertirse en servicios de datos a prueba de fallos o escalables.

RGM también gestiona los recursos del clúster, incluidas las instancias de una aplicación y de los recursos de red (nombres de sistema lógicos y direcciones compartidas).

Además de los métodos proporcionados por el software Sun Cluster, el sistema SunPlex también proporciona una API y varias herramientas de desarrollo de servicios de datos que permiten a los programadores de aplicaciones desarrollar los métodos de servicio de datos necesarios para que otras aplicaciones se ejecuten como servicios de datos de alta disponibilidad con el software de Sun Cluster.

Servicios de datos a prueba de fallos

Si el nodo en el que se está ejecutando el servicio de datos (el primario) falla, el servicio migra a otro nodo en funcionamiento sin intervención del usuario. Los servicios a prueba de fallos usan un *grupo de recursos de recuperación de fallos*, que es un contenedor para recursos de instancias de aplicaciones y de red (*nombres de sistema lógicos*). Éstos son direcciones IP que pueden configurarse como activas en un nodo y, posteriormente, configurarse automáticamente como inactivas en el nodo original y activarse en otro nodo.

En los servicios de datos a prueba de fallos, las instancias de las aplicaciones sólo se ejecutan en un nodo individual. Si el supervisor de fallos detecta un error, intenta reiniciar la instancia del mismo nodo o la de otro nodo (recuperación de fallos), dependiendo de la forma en que se haya configurado el servicio de datos.

Servicios de datos escalables

Los servicios de datos escalables tienen la capacidad de activar instancias en varios nodos; usan dos grupos de recursos: un *grupo de recursos escalable* que contiene los recursos de aplicaciones y un grupo de recursos a prueba de fallos que contiene los recursos de red (*direcciones compartidas*) de los que depende un servicio escalable. El grupo de recursos escalable puede estar en línea en varios nodos, de forma que se puedan ejecutar varias instancias del servicio simultáneamente. El grupo de recurso a prueba de fallos que aloja la dirección compartida está disponible en un solo nodo cada vez. Todos los nodos que alojan un servicio escalable usan la misma dirección compartida para alojar el servicio.

Las peticiones de servicio llegan al clúster a través de una interfaz de red individual (interfaz global) y se distribuyen a los nodos según uno de varios algoritmos definidos por la *política de equilibrio de cargas* que el clúster puede usar para equilibrar la carga del servicio entre varios nodos. Tenga en cuenta que pueden haber varias interfaces globales en distintos nodos alojando otras direcciones compartidas.

En los servicios escalables, las instancias de las aplicaciones se ejecutan en varios nodos simultáneamente. Si el nodo que aloja la interfaz global falla, ésta se traspassa a otro nodo. Si falla una instancia de aplicación en ejecución, ésta intenta reiniciarse en el mismo nodo.

Si no se puede reiniciar una instancia de una aplicación en el mismo nodo, y se configura otro nodo que no esté en uso para ejecutar el servicio, éste se transfiere a este nodo. En caso contrario, continúa ejecutándose en los nodos restantes, lo que podría provocar una degradación en el rendimiento del servicio.

Nota – El estado TCP para cada instancia de aplicación se conserva en el nodo de la instancia, no en el de la interfaz global. Por tanto el fallo en el nodo de interfaz global no afecta a la conexión.

En la [Figura 3-7](#) se muestra un ejemplo de grupo de recursos escalable y las dependencias que existen entre ellos para los servicios escalables. Este ejemplo muestra tres grupos de recursos. El grupo de recursos a prueba de fallos contiene recursos de aplicación que usan los DNS de alta disponibilidad y recursos de red que usan éstos y el servidor Apache Web Server de alta disponibilidad. Los grupos de recursos escalables sólo contienen instancias de la aplicación de Apache Web Server. Tenga en cuenta que existen dependencias de grupos de recursos entre los grupos de recursos escalables y los a prueba de fallos (líneas sólidas) y que todos los demás recursos de la aplicación de Apache dependen del recurso de red `schost-2`, que es una dirección compartida (líneas punteadas).

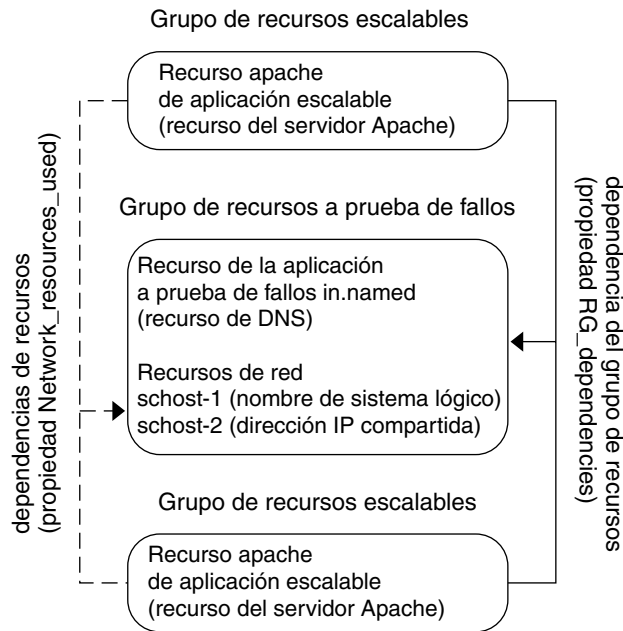


FIGURA 3-7 SPARC: Ejemplo de los grupos de recursos a prueba de fallos y escalables

Política de equilibrio de cargas

El equilibrio de cargas mejora el rendimiento del servicio escalable, tanto en tiempo de respuesta como en rendimiento.

Existen dos clases de servicios de datos escalables: *puros* y *adosados*. Un servicio puro es aquel cuyas instancias pueden responder a peticiones de clientes sin restricciones. Un servicio adosado es uno en el que un cliente envía peticiones a la misma instancia. Estas peticiones no se redirigen a otras instancias.

Un servicio puro usa una política de equilibrio de cargas ponderada bajo la cual, predeterminadamente, las peticiones de los clientes se distribuyen de manera uniforme entre las instancias del servidor en el clúster. Por ejemplo en un cluster de tres nodos supongamos que cada nodo pesa una unidad. Cada nodo dará servicio a 1/3 de las peticiones de cualquier cliente en nombre de ese servicio. El administrador puede cambiar los pesos en todo momento con la interfaz de órdenes `scrgadm (1M)` o con la GUI de administración de SunPlex.

Existen dos variedades de servicios adosados, *normales* y *comodín*, ambos permiten que sesiones simultáneas de aplicación a través de varias conexiones TCP compartan el estado de la memoria (estado de la sesión de aplicación).

Los normales permiten a un cliente compartir el estado entre varias conexiones TCP simultáneas. Al cliente se le denomina "adosado" con respecto a la instancia del servidor que está a la escucha en un único puerto. Al cliente se le garantiza que todas sus solicitudes vayan a la misma instancia del servidor, siempre que ésta permanezca activa y accesible y que la política de equilibrio de cargas no cambie mientras el servicio esté en línea.

Por ejemplo, un explorador web en el cliente se conecta a una dirección IP compartida en el puerto 80 usando tres conexiones TCP distintas, pero las conexiones están intercambiando información de sesión en memoria intermedia entre ellas en el servicio.

Una generalización de política de adosamiento se aplica a varios servicios escalables que intercambien información de forma no visible en la misma instancia. Cuando estos servicios intercambian información de sesión de forma no visible en la misma instancia, se dice que el cliente está "adosado" con respecto a varias instancias de servidor del mismo nodo que está a la escucha en puertos distintos.

Por ejemplo, un cliente de una sede de comercio electrónico rellena su carro de la compra con artículos usando HTTP normal en el puerto 80, pero cambia a SSL en el puerto 443 para enviar datos seguros y pagar por tarjeta de crédito los artículos del carro de la compra.

Los servicios adosados comodín usan números de puerto asignados dinámicamente, pero siguen esperando que las peticiones de clientes vayan al mismo nodo. El cliente está "adosado con comodín" a los puertos con respecto a la misma dirección IP.

Un buen ejemplo de esta política es el FTP de modalidad pasiva. Un cliente se conecta a un servidor FTP en el puerto 21 y después el servidor le informa que vuelva a conectarse a un servidor de puertos que está a la escucha en el rango de puertos dinámico. Todas las solicitudes para esta dirección IP se redirigen al mismo nodo que el servidor informó al cliente a través de la información de control.

Tenga en cuenta que además de estas políticas de adosamiento también está vigente de forma predeterminada la política de equilibrio de cargas ponderado, por tanto la petición inicial de un cliente se dirige a la instancia que dicta el repartidor de cargas. Cuando el cliente ha establecido una afinidad para el nodo en que se está ejecutando la instancia, las futuras peticiones se dirigen a esa instancia siempre que el nodo esté accesible y la política de equilibrio de cargas no cambie.

A continuación se explican detalles adicionales sobre las políticas de equilibrio de cargas específicas.

- **Ponderada:** la carga se distribuye entre varios nodos según valores de peso especificados. Esta política se establece mediante el valor `LB_WEIGHTED` para la propiedad `Load_balancing_weights`. Si no se establece explícitamente el peso de un nodo, éste toma de forma predeterminada el valor uno.

La política de ponderación redirige cierta parte del tráfico de los clientes a un nodo determinado. Dado X =peso y A =peso total de todos los nodos activos, cada uno puede esperar recibir aproximadamente X/A del total de conexiones nuevas, cuando éste sea lo bastante grande. Esta política no trata de peticiones individuales.

Tenga en cuenta que esta política no es de tipo giratoria. Una política de puerta giratoria siempre haría que todas las peticiones de un cliente fueran a un nodo distinto: la primera petición al nodo 1, la segunda al nodo 2, etc.

- **Adosada:** en esta norma el conjunto de puertos se da a conocer en el momento en el que se configuran los recursos de la aplicación. Se establece mediante el valor `LB_STICKY` para la propiedad del recurso `Load_balancing_policy`.
- **Adosada-comodín:** esta política tiene menos limitaciones que la “adosada” normal. En un servicio escalable identificado por su dirección IP, el servidor asigna los puertos (y no se conocen de antemano). Los puertos podrían cambiar. Esta política se establece mediante el valor `LB_STICKY_WILD` para la propiedad del recurso `Load_balancing_policy`.

Valores de retroceso

Los grupos de recurso se cubren entre nodos. Cuando esto se produce, el que hacía el papel de secundario se convierte en el nuevo primario. La configuración de retroceso especifica las acciones que acontecen cuando el primario original vuelve a estar en línea. Las opciones son hacer que el primario original se convierta de nuevo en primario (retroceso) o permitir que el que haya tomado su papel continúe con él. La opción deseada se especifica mediante el valor `Failback` de la propiedad del grupo de recurso.

En algunos casos si el nodo original que aloja al grupo de recursos está fallando y rearrancando repetidas veces, configurar el retroceso podría producir una disponibilidad menor para el grupo de recursos.

Supervisores de fallos de servicios de datos

Cada servicio de datos de SunPlex proporciona un supervisor de fallos que explora periódicamente el servicio de datos para determinar su buen estado. Un supervisor de fallos comprueba que los daemons de las aplicaciones estén funcionando y que se esté dando servicio a los clientes. De acuerdo con la información que devuelven las sondas, se pueden tomar acciones predefinidas, como reiniciar daemons o producir una recuperación de fallos.

Desarrollo de nuevos servicios de datos

Sun proporciona archivos de configuración y plantillas de métodos de gestión que permiten hacer funcionar varias aplicaciones como servicios a prueba de fallos o escalables dentro de un clúster. Si Sun no ofrece en algún momento una aplicación que se desee ejecutar como servicio a prueba de fallos o escalable, se puede usar la API DSET u otra para configurar la aplicación a fin de que se ejecute como un servicio de tipo a prueba de fallos o escalable.

Existe un conjunto de criterios para determinar si una aplicación puede convertirse en un servicio a prueba de fallos. Los criterios específicos se describen en los documentos de SunPlex que tratan de las API que se puede usar para las aplicaciones.

Aquí, presentamos algunas directrices que permitirán comprender si el servicio puede aprovechar las ventajas de la arquitectura de servicios de datos escalable. Repase el apartado [“Servicios de datos escalables” en la página 68](#) para obtener más información sobre los servicios escalables.

Los servicios nuevos que sigan estas directrices pueden utilizar las ventajas de los servicios escalables. Si un servicio existente no sigue estas directrices exactamente, será necesario reescribir partes del mismo para que las cumpla.

Un servicio de datos escalable tiene las características siguientes. En primer lugar, un servicio así está compuesto por una o más *instancias* de servidor, cada una de las cuales se ejecuta en un nodo distinto del clúster. Dos o más instancias del mismo servicio no pueden ejecutarse en el mismo nodo.

En segundo lugar, si el servicio ofrece un almacenamiento de datos lógico externo, debe sincronizarse el acceso simultáneo a este almacenamiento desde varias instancias de servidores para evitar perder actualizaciones o leer los datos mientras se estén

modificando. Tenga en cuenta que se llama “externo” para distinguir el almacenamiento de estado en memoria y “lógico” porque parece como una entidad única, aunque puede estar replicado. Además, este almacenamiento de datos lógico tiene la propiedad de que cuando alguna instancia de servidor lo actualiza, las demás instancias pueden ver las modificaciones .

El sistema SunPlex proporciona ese almacenamiento externo a través de su sistema de archivos del clúster y las particiones a bajo nivel globales. Como ejemplo, supongamos que un servicio escribe datos nuevos a un archivo de registro cronológico externo o modifica los datos que haya. Cuando se ejecuten varias instancias de este servicio, cada una tendrá acceso a este registro cronológico externo y podrá intentar acceder al mismo simultáneamente. Cada instancia debe sincronizar su acceso a este registro o de lo contrario se interferirán entre ellas. El servicio podría usar bloqueos de archivo normales de Solaris con `fcntl(2)` y `lockf(3C)` para conseguir la sincronización deseada.

Otro ejemplo de este tipo de almacenamiento es la base de datos de servidores, como Oracle de alta disponibilidad o Oracle Real Application Clusters para clústeres basados en SPARC. Tenga en cuenta que este tipo de servidor de base de datos de componente trasero proporciona una sincronización incorporada gracias a las consultas de base de datos o las transacciones de actualización, por lo que no es necesario que las distintas instancias de servidor implementen su propia sincronización.

Un ejemplo de servicio que no es de tipo escalable en su encarnación actual es el servidor IMAP de Sun. El servicio actualiza un almacenamiento, pero es privado y cuando varias instancias de IMAP escriben en él, se sobrescriben porque las actualizaciones no están sincronizadas. El servidor IMAP debe volver a escribirse para que se sincronice el acceso simultáneo.

Finalmente, tenga en cuenta que algunas instancias pueden tener datos privados que sean contradictorios con los de otras. En tal caso, no es necesario que el servicio se preocupe de sincronizar el acceso simultáneo, ya que los datos son privados y sólo esa instancia puede manipularlos. En este caso, se ha de tener cuidado de no almacenar estos datos privados en el sistema de archivos del clúster, porque existe la posibilidad de que resulte accesible globalmente.

API de servicio de datos y de biblioteca de desarrollo de servicio de datos

El sistema SunPlex para que las aplicaciones estén altamente disponibles ofrece:

- Servicios de datos proporcionados como parte del sistema SunPlex
- Una API de servicio de datos
- Una API de biblioteca de desarrollo de servicio de datos
- Un servicio de datos “genérico”

Sun Cluster Data Services Planning and Administration Guide for Solaris OS describe cómo instalar y configurar los servicios de datos proporcionados por el sistema SunPlex. *Sun Cluster 3.1 9/04 Software Collection for Solaris OS (SPARC Platform Edition)* describe cómo configurar otras aplicaciones para que tengan una alta disponibilidad en la arquitectura Sun Cluster.

Las API de Sun Cluster permiten a los programadores de aplicaciones desarrollar supervisores de fallos y secuencias que inician y detienen instancias de servicios de datos. Con estas herramientas una aplicación puede instrumentalizarse para que sea un servicio a prueba de fallos o escalable. Además, el sistema SunPlex ofrece un servicio de datos "genérico" que se puede usar para generar rápidamente métodos de inicio y detención necesarios para la aplicación y ejecutarla como servicio a prueba de fallos o escalable.

Uso de la interconexión del clúster para el tráfico de servicio de datos

Un clúster debe tener varias conexiones de red entre los nodos que formen la interconexión del clúster. El software de agrupación en clúster usa varias interconexiones para la alta disponibilidad y para mejorar el rendimiento. Para el tráfico interno (por ejemplo, datos del sistema de archivos o datos de servicios escalables), los mensajes se reparten entre todas las interconexiones disponibles como si pasaran a través de una puerta giratoria.

La interconexión del clúster también está disponible para aplicaciones, para comunicaciones de alta disponibilidad entre nodos. Por ejemplo, una aplicación distribuida podría tener componentes que se ejecutaran en distintos nodos y que necesiten comunicarse. Al usar la interconexión del clúster en lugar del transporte público, éstas conexiones pueden soportar el fallo de un enlace individual.

Para usar la interconexión del clúster para la comunicación entre nodos, una aplicación necesita los nombres de sistema privados configurados cuando se instaló el clúster. Por ejemplo, si el nombre de sistema privado para el nodo 1 es `clusternode1-priv`, use ese nombre para comunicarse a través de la interconexión del clúster con el nodo 1. Los zócalos TCP abiertos mediante este nombre se dirigen por la interconexión del clúster y pueden re-dirigirse de forma transparente en caso de un fallo de la red

Tenga en cuenta que, debido a que los nombres de sistema privados pueden configurarse durante la instalación, la interconexión del clúster podría usar cualquiera de los nombres elegidos en ese momento. El nombre real podrá obtenerse con `scha_cluster_get(3HA)` con el argumento `scha_privatelink_hostname_node`.

Para el uso de la interconexión del clúster a nivel de aplicación, se usa una interconexión simple entre cada par de nodos, aunque si es posible es mejor utilizar interconexiones distintas para distintos pares de nodos. Por ejemplo, consideremos una aplicación que se ejecuta en tres nodos basados en plataformas SPARC y se comunica mediante la interconexión del clúster. La comunicación entre los nodos 1 y 2 podría darse en la interfaz `hme0`, mientras que la comunicación entre los nodos 1 y 3 podría producirse en la interfaz `qfe1`. Es decir, que la comunicación entre dos nodos estaría limitada a la interconexión simple, mientras que internamente la comunicación del clúster se dividiría entre todas las interconexiones.

Tenga en cuenta que la aplicación comparte la interconexión con el tráfico interno del clúster, por lo que el ancho de banda disponible para la aplicación depende del que se use para el resto del tráfico del clúster. En caso de fallo, el tráfico interno puede desviarse por el resto de interconexiones, mientras que las conexiones de aplicación de una interconexión fallida se cambian a otra que funcione.

Hay dos tipos de direcciones que admiten la interconexión del clúster y `gethostbyname(3N)` en un nombre de sistema privado normalmente devuelve dos direcciones IP. La primera dirección se denomina *dirección pairwise lógica* y la segunda *dirección pernode lógica*.

A cada par de nodos se le asigna una dirección lógica pairwise independiente. Esta pequeña red lógica admite la recuperación de fallos de las conexiones. A cada nodo también se le asigna una dirección pernode fija. Es decir, las direcciones pairwise lógicas para `clusternode1-priv` son distintas para cada nodo, mientras que la dirección pernode lógica para `clusternode1-priv` es la misma en todos los nodos. Un nodo no dispone de dirección pairwise consigo mismo, por tanto `gethostbyname(clusternode1-priv)` en el nodo 1 sólo devuelve la dirección pernode lógica.

Tenga en cuenta que las aplicaciones que acepten conexiones a través de la interconexión del clúster y que después verifiquen la dirección IP por motivos de seguridad deben comprobarse con todas las direcciones IP que haya devuelto la orden `gethostbyname` y no sólo con la primera de ellas.

Si necesita direcciones IP uniformes en su aplicación en todos los puntos, configure la aplicación para vincularla con la dirección pernode en los lados del cliente y el servidor para que parezca que todas las conexiones vengan y vayan de la dirección pernode.

Recursos, grupos de recursos y tipos de recursos

Los servicios de datos usan varios tipos de *recursos*: las aplicaciones como Sun Java System Web Server (anteriormente Sun Java System Web Server) o Apache Web Server usan direcciones de red (nombres de sistema lógicos y direcciones compartidas) de las que dependen las aplicaciones. Los recursos de aplicación y red forman una unidad básica que gestiona RGM.

Los servicios de datos son tipos de recursos. Por ejemplo, Sun Cluster HA for Oracle es el tipo de recurso `SUNW.oracle-server` y Sun Cluster HA for Apache es `SUNW.apache`.

Nota – El tipo de recurso `SUNW.oracle-server` sólo se usa en clústers basados en plataformas SPARC.

Un recurso es una concreción de *tipo de recurso* que está definida a nivel del clúster. Hay definidos distintos tipos de recursos.

Los recursos de red son tipos de recurso `SUNW.LogicalHostname` o `SUNW.SharedAddress`. Ambos están registrados previamente por el software de Sun Cluster.

Los tipos de recurso `SUNW.HAStorage` y `HAStoragePlus` se usan para sincronizar el inicio de recursos y los grupos de dispositivos de disco de los que éstos dependen. Aseguran que antes de que se inicie un servicio de datos, estén disponibles las rutas de acceso a los puntos de montaje de los sistemas de archivos del clúster, dispositivos globales y grupos de dispositivos. Para obtener más información, consulte “Synchronizing the Startups Between Resource Groups and Disk Device Groups” en el manual *Data Services Installation and Configuration Guide*. (El tipo de recurso `HAStoragePlus` estuvo disponible en Sun Cluster 3.0 5/02 e incorporó otra característica, permitiendo a los sistemas de archivo locales que estuvieran altamente disponibles. Para obtener más información sobre esta función, consulte “[Tipo de recurso HAStoragePlus](#)” en la página 49.).

Los recursos gestionados por RGM se sitúan en grupos denominados *grupos de recursos*, de forma que pueden ser gestionados como una unidad. El grupo de recursos migra como unidad si se inicia una recuperación de fallos o un cambio.

Nota – Cuando un grupo de recursos que contiene recursos de aplicación se pone en línea, la aplicación se inicia. El método de inicio del servicio de datos espera hasta que la aplicación esté completamente en marcha antes de salir con éxito. Determinar cuándo la aplicación está completamente en marcha sigue el mismo proceso que el supervisor de fallos utiliza para saber que un servicio de datos está ofreciendo servicio a clientes. Consulte *Sun Cluster Data Services Planning and Administration Guide for Solaris OS* para obtener más información sobre este proceso.

Resource Group Manager (RGM)

RGM controla los servicios de datos (aplicaciones) como recursos, que las implementaciones de *tipos de recursos* gestionan. Éstas las proporciona Sun o las crea un desarrollador con una plantilla genérica de servicios de datos, una API de la biblioteca de desarrollo de servicios de datos (API DSDL) o con una API de gestión de recursos (RMAPI). El administrador del clúster crea y gestiona recursos en contenedores llamados *grupos de recursos*. RGM para e inicia los grupos de recursos en nodos seleccionados como respuesta a cambios en la pertenencia al clúster.

RGM actúa en *recursos* y *grupos de recursos*. Las acciones de RGM hacen que los recursos y los grupos de recursos cambien entre los estados en línea y fuera de línea. En el apartado [“Estados y configuración de recursos y grupos de recursos” en la página 77](#) puede encontrarse una descripción completa de los estados y valores que pueden aplicarse a recursos y grupos de recursos. Consulte [“Recursos, grupos de recursos y tipos de recursos” en la página 76](#) para obtener información sobre cómo ejecutar un proyecto de gestión de recursos bajo el control de RGM.

Estados y configuración de recursos y grupos de recursos

Los administradores aplican a recursos y grupos de recursos configuraciones estáticas que sólo pueden cambiarse con acciones administrativas. RGM cambia los grupos de recursos entre los estados “dinámicos.” Estos valores y estados se describen en la lista siguiente.

- **Gestionados o no gestionados:** son valores que afectan a todo el clúster y sólo se aplican a grupos de recursos. Los grupos de recursos los gestiona RGM. La orden `scrgadm (1M)` se puede utilizar para provocar que RGM se encargue o no de la gestión de un grupo de recursos. Estos valores no cambian con la reconfiguración del clúster.

Cuando se crea un grupo de recursos por primera vez, no se gestiona. Debe gestionarse antes de que cualquier recurso que se incluya en el grupo pueda activarse.

En algunos servicios de datos, por ejemplo en un servidor web escalable, el trabajo debe hacerse antes de iniciar los recursos de red y después de que se detengan. Este trabajo se hace con los métodos de servicio de datos de inicialización (INIT) y finalización (FINI). Los métodos INIT sólo se ejecutan si el grupo de recursos en el que éstos residen está en estado gestionable.

Cuando un grupo de recursos cambia de no gestionado a gestionado, los métodos INIT registrados para el grupo se ejecutan en todos los recursos.

Cuando un grupo de recursos cambia de gestionado a no gestionado, todos los métodos FINI registrados se ejecutan para realizar una limpieza.

El uso más común de los métodos INIT y FINI son para recursos de red de servicios escalables, pero pueden usarse para cualquier trabajo de inicialización o limpieza que no realice la propia aplicación.

- **Habilitados o inhabilitados:** son los valores a nivel del clúster que se aplican a los recursos. El comando `scrgadm (1M)` se puede usar para habilitar o inhabilitar los recursos. Estos valores no cambian con la reconfiguración del clúster.

El valor normal para un recurso es que esté habilitado y funcionando activamente en el sistema.

Si por algún motivo desea que el recurso no esté disponible en todos los nodos del clúster, puede inhabilitarlo. De esta manera dejará de estar disponible para uso general.

- **En línea o fuera de línea:** son estados dinámicos y se aplican a recursos y grupos de recursos.

Estos estados cambian a medida que el clúster pasa a través de los pasos de reconfiguración durante una conmutación o una recuperación de fallos. También pueden cambiarse a través de acciones de administración. La orden `scswitch (1M)` se puede usar para cambiar los estados en línea o fuera de línea de un recurso o grupo de recursos.

Un recurso o un grupo de recursos a prueba de fallos sólo pueden estar en línea en un nodo simultáneamente. Un recurso o un grupo de recursos escalables pueden estar en línea en algunos nodos y fuera de línea en otros. Durante una conmutación o una recuperación de fallos, los grupos de recursos y los recursos que contienen se ponen en fuera de línea en un nodo y después se vuelven a poner en línea en otro distinto.

Si un grupo de recursos está fuera de línea significa que todos sus recursos también lo están. Si un grupo de recursos está en línea significa que todos sus recursos habilitados también lo están.

Los grupos de recursos pueden contener varios recursos, pudiendo haber dependencias entre ellos que requieren que los recursos se pongan en línea y fuera de línea en un orden determinado. Los métodos usados para poner los recursos en línea y fuera de línea pueden ocupar tiempos distintos en cada uno de ellos. Debido a las dependencias de recursos y a las diferencias de tiempo de inicio y finalización, los recursos de un mismo grupo de recursos pueden tener estados de puesta en línea y fuera de línea distintos durante una reconfiguración del clúster.

Propiedades de recursos y grupos de recursos

En los servicios de datos de SunPlex se pueden configurar valores de propiedad para recursos y grupos de recursos. Las propiedades estándar son comunes a todos los servicios de datos. Las propiedades de extensión son específicas de cada servicio de datos. Algunas propiedades estándar y de extensión se configuran con valores predeterminados para que no se tengan que modificar. Otras necesitan configurarse como parte del proceso de creación y configuración de recursos. La documentación de cada servicio de datos especifica qué propiedades de recurso pueden establecerse y cómo hacerlo.

Las propiedades estándar se usan para configurar propiedades de recursos y grupos de recursos que normalmente son independientes de todos los servicios de datos. Para obtener información sobre el conjunto de propiedades estándar, consulte “Standard Properties” in *Sun Cluster Data Services Planning and Administration Guide for Solaris OS*.

Las propiedades de extensión de RGM ofrecen información como la ubicación de los binarios de la aplicación y los archivos de configuración. Las propiedades de extensión se modifican a medida que se configuran los servicios de datos. El conjunto de propiedades de extensión se describe en la guía específica para servicios de datos.

Configuración del proyecto de servicios de datos

Los servicios de datos pueden configurarse para que se ejecuten bajo un nombre de proyecto de Solaris cuando se pongan en línea con RGM. La configuración asocia un recurso o un grupo de recursos gestionados por RGM con un OD de proyecto de Solaris. La correlación desde el recurso o grupo de recursos a un ID de proyecto ofrece la posibilidad de usar controles sofisticados que están disponibles en el entorno Solaris para gestionar las cargas de trabajo y el consumo dentro del clúster.

Nota – Esta configuración sólo se puede realizar si se está ejecutando la versión actual del software de Sun Cluster con Solaris 9.

La funcionalidad de la gestión de Solaris en un entorno clúster permite dar prioridad a las aplicaciones más importantes al compartir un nodo con otras aplicaciones. Las aplicaciones podrían compartir un nodo si hay servicios consolidados o porque se haya producido una recuperación de fallos. El uso de la funcionalidad de la gestión que se describe aquí podría mejorar la disponibilidad de una aplicación crítica ya que evita que otras aplicaciones de prioridad baja consuman demasiados suministros del sistema, como tiempo de CPU.

Nota – La documentación de Solaris sobre esta función describe el tiempo de CPU, los procesos, las tareas y componentes similares como 'recursos'. Sin embargo, la documentación de Sun Cluster usa el término 'recursos' para describir entidades que están bajo el control de RGM. El apartado siguiente usará el término 'recurso' para hacer referencia a las entidades de Sun Cluster bajo el control de RGM y el término 'suministros' para hacer referencia a tiempo de CPU, procesos y tareas.

Este apartado ofrece una descripción conceptual de la configuración de los servicios de datos para ejecutar procesos en un `project(4)` especificado de Solaris 9. Este apartado también describe varios escenarios de recuperación de fallos y sugerencias para planificar el uso de la funcionalidad de la gestión que ofrece el entorno Solaris. Para obtener documentación detallada sobre conceptos y procedimientos de la función de la gestión, consulte *System Administration Guide: Resource Management and Network Services* en *Solaris 9 System Administrator Collection*.

Al configurar recursos y grupos de recursos para usar la funcionalidad de la gestión de Solaris en un clúster, considere el uso del siguiente proceso de alto nivel:

1. Configure aplicaciones como parte del recurso.
2. Configure recursos como parte de un grupo de recursos.
3. Habilite los recursos del grupo.
4. Haga gestionable el grupo de recursos.
5. Cree un proyecto de Solaris para el grupo de recursos.
6. Configure propiedades estándar para asociar el nombre del grupo de recursos con el proyecto que se ha creado en el paso 5.
7. Ponga en línea el grupo de recursos.

Para configurar las propiedades estándar `Resource_project_name` o `RG_project_name` para asociar el ID de proyecto de Solaris con el recurso o grupo de recursos, use la opción `-y` con la orden `scrgradm(1M)`. Configure los valores de la propiedad al recurso o grupo de recursos. Consulte "Standard Properties" en *Sun Cluster Data Services Planning and Administration Guide for Solaris OS* para obtener las descripciones adecuadas. Consulte `r_properties(5)` y `rg_properties(5)` para obtener las descripciones adecuadas.

El nombre del proyecto especificado debe existir en la base de datos de proyectos (`/etc/project`) y el superusuario debe estar configurado como miembro del proyecto con ese nombre. Consulte "Projects and Tasks" en *System Administration Guide: Resource Management and Network Services* en *Solaris 9 System Administrator Collection* para obtener información conceptual sobre el nombre del proyecto. Consulte `project(4)` para obtener una descripción de la sintaxis de los archivos del proyecto.

Cuando RGM pone en línea el recurso o grupo de recursos, ejecuta los procesos relacionados que hay bajo el nombre del proyecto.

Nota – Los usuarios pueden asociar recursos o grupos de recursos con proyectos en cualquier momento. Sin embargo, el nombre del proyecto nuevo no tiene vigencia hasta que el recurso o grupo de recursos se pone fuera de línea y se vuelve a poner en línea usando RGM.

Ejecutar recursos y grupos de recursos bajo el nombre del proyecto permite configurar las funciones siguientes para gestionar suministros de sistema en todo el clúster.

- Contabilidad ampliada: ofrece una forma flexible de registrar el consumo según las tareas o los procesos. La contabilidad ampliada permite examinar el uso histórico y evaluar las necesidades de capacidad para cargas de trabajo futuras.
- Controles: proporcionan un mecanismo para ajustarse a los suministros del sistema. Puede evitarse que los procesos, tareas y proyectos consuman grandes cantidades de suministros de sistema especificados.
- Programación de partición justa (FSS): ofrece la posibilidad de controlar la ubicación de tiempo de CPU disponible entre cargas de trabajo según su importancia. Ésta se mide por el número de particiones de tiempo de CPU que se asigna a cada carga de trabajo. Consulte `dispadm(1M)` para obtener una descripción de las líneas de órdenes para establecer FSS como programador predeterminado. Consulte también `priocntl(1)`, `ps(1)` y `FSS(7)` para obtener más información.
- Agrupaciones: ofrecen la posibilidad de usar particiones para aplicaciones interactivas según los requisitos de éstas. La agrupaciones pueden usarse para particionar un servidor que admite distintas aplicaciones de software. El uso de agrupaciones produce una respuesta más predecible para cada aplicación.

Determinación de requisitos para la configuración de proyectos

Antes de que se configuren servicios de datos para usar los controles que ofrece Solaris en un entorno Sun Cluster, es necesario decidir cómo se desean controlar y seguir los recursos en caso de cambios o recuperación de fallos. Considere identificar dependencias dentro del clúster antes de configurar un proyecto nuevo. Por ejemplo, los recursos y grupos de recursos dependen de grupos de dispositivos de disco. Las propiedades del grupo de recursos `nodelist`, `failback`, `maximum primaries` y `desired primaries` se configuran con `scrgadm(1M)` para identificar prioridades de las listas de nodos en el grupo de recursos. Consulte “Relationship Between Resource Groups and Disk Device Groups” en *Sun Cluster Data Services Planning and Administration Guide for Solaris OS* para obtener una explicación breve sobre las dependencias de la lista de nodos entre los grupos de recursos y los grupos de dispositivos del disco. Para obtener descripciones de propiedad detalladas, consulte `rg_properties(5)`.

Utilice las propiedades `preferred` y `failback` configuradas con `scrgadm(1M)` y `scsetup(1M)` para determinar las prioridades de la lista de nodos del grupo de dispositivos de disco. Para obtener información sobre procedimientos, consulte “Cómo cambiar las propiedades de un dispositivo de disco” en “Administering Disk Device Groups” en *Sun Cluster System Administration Guide for Solaris OS*. Consulte “Los componentes del hardware y el software de SunPlex” en la página 21 para obtener información conceptual sobre la configuración de los nodos y el comportamiento de los servicios de datos a prueba de fallos y escalables.

Si se configura todos los nodos del clúster de forma idéntica, los límites de uso se hacen cumplir de forma idéntica en los nodos primarios y secundarios. Es necesario que los parámetros de configuración de los proyectos no sean idénticos en todas las aplicaciones de todos los nodos. Todos los proyectos asociados con la aplicación deben ser accesibles al menos por la base de datos de proyectos en todos los maestros potenciales de esa aplicación. Supongamos que la aplicación 1 está controlada por `phys-schost-1` pero podría efectuarse una conmutación o una resolución de fallo a `phys-schost-2` o `phys-schost-3`. El proyecto asociado con la aplicación 1 debe estar accesible en los tres nodos (`phys-schost-1`, `phys-schost-2` y `phys-schost-3`).

Nota – La información de la base de datos de proyectos puede ser un archivo local `/etc/project` o puede estar almacenada en el mapa NIS o en el servicio de directorios LDAP.

El entorno Solaris permite una configuración flexible de parámetros de uso y Sun Cluster pone pocas restricciones. Las opciones de configuración dependen de las necesidades de la sede. Considere las directrices generales de los apartados siguientes antes de configurar los sistemas.

Establecimiento de límites de memoria virtual por proceso

Configure el control `process.max-address-space` para limitar la memoria virtual a cada proceso. Consulte `rctladm(1M)` para obtener información detallada sobre la configuración del valor `process.max-address-space`.

Al utilizar controles de gestión con Sun Cluster, se deben configurar los límites de memoria apropiadamente para evitar una recuperación de fallos innecesaria de aplicaciones y un efecto “ping-pong”. En general:

- No configure unos límites de memoria demasiado bajos.
 - Cuando una aplicación llegue a su límite de memoria, podría producirse una recuperación de fallos. Esta directriz es especialmente importante para aplicaciones de base de datos, ya que al alcanzar un límite de memoria virtual se pueden producir consecuencias inesperadas.

- No configure límites de memoria idénticos en nodos primarios y secundarios.
Límites idénticos podrían causar un efecto ping-pong cuando una aplicación alcance su límite de memoria y se recupere el error en un nodo secundario con un límite de memoria idéntico. Configure el límite de memoria un poco por encima en el nodo secundario. La diferencia en los límites de memoria ayuda a evitar la situación ping-pong y da al administrador del sistema un periodo de tiempo en el que ajustar los parámetros como sea necesario.
- Use los límites de memoria de la gestión de recursos para el equilibrio de cargas.
Por ejemplo, puede usar límites de memoria para evitar que una aplicación que no funciona bien consuma demasiado espacio de intercambio en disco.

Casos de recuperación de fallos

Puede configurar parámetros de gestión para que la asignación en la configuración del proyecto (`/etc/project`) funcione normalmente en el clúster y en las situaciones de conmutación o recuperación de fallos.

Los apartados siguientes son casos de ejemplo.

- Los dos primeros apartados, “Clúster de dos nodos con dos aplicaciones” y “Clúster de dos nodos con tres aplicaciones”, muestran escenarios de recuperación de fallos para nodos completos.
- El apartado “Recuperación de fallos sólo de grupos de recursos” ilustra el funcionamiento de la recuperación de fallos sólo para una aplicación.

En un entorno del clúster las aplicaciones se configuran como parte de un recurso y éste como parte de un grupo de recursos (RG). Cuando se produce un fallo, el grupo de recursos, junto con su aplicación asociada, se transfieren a otro nodo. En los ejemplos siguientes los recursos no se muestran específicamente. Se presupone que cada recurso sólo tiene una aplicación.

Nota – La recuperación de fallos se produce en el orden de lista de nodos especificado en la RGM.

Los ejemplos siguientes tienen estas limitaciones:

- La aplicación 1 (App-1) está configurada en el grupo de recursos RG-1.
- La aplicación 2 (App-2) está configurada en el grupo de recursos RG-2.
- La aplicación 3 (App-3) está configurada en el grupo de recursos RG-3.

Aunque los números de particiones asignadas siguen siendo los mismos, el porcentaje de tiempo de CPU asignado a cada aplicación cambia después de la recuperación de fallos. Este porcentaje depende del número de aplicaciones que se están ejecutando en el nodo y del número de particiones que se han asignado a cada aplicación activa.

En estos casos, se asumen las configuraciones siguientes.

- Todas las aplicaciones están configuradas bajo un proyecto común.
- Cada recurso dispone de una única aplicación.
- Las aplicaciones son los únicos procesos activos de los nodos.
- Las bases de datos de los proyectos están configuradas igual en todos los nodos del clúster.

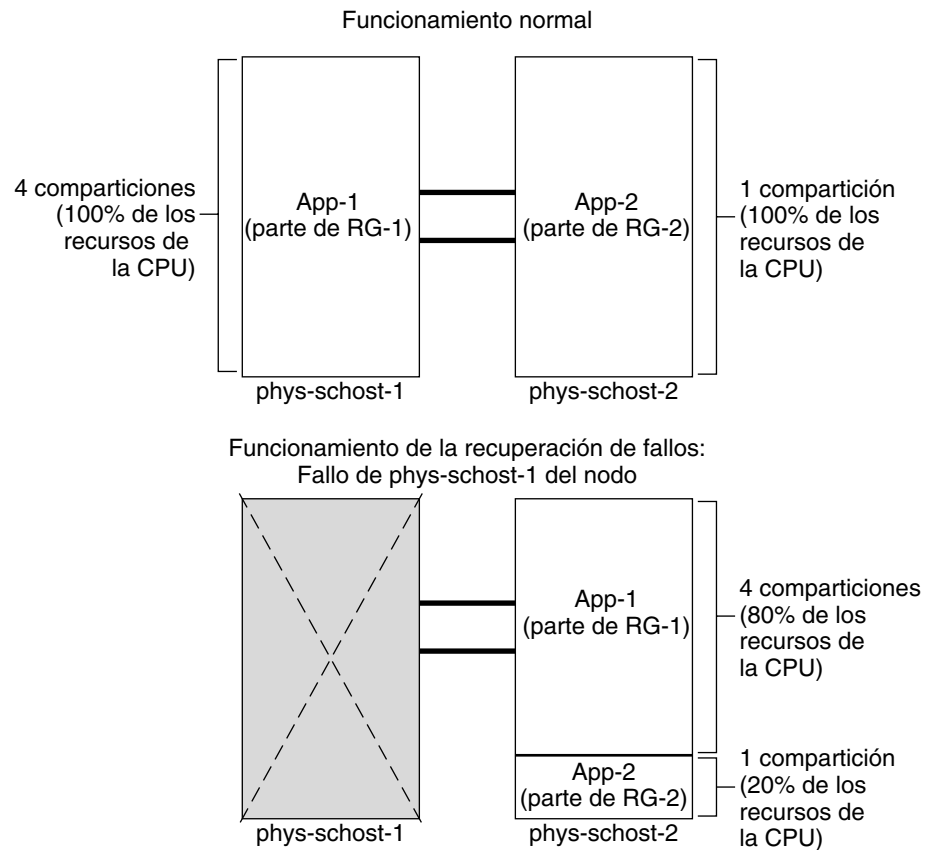
Clúster de dos nodos con dos aplicaciones

Se pueden configurar dos aplicaciones en un clúster de dos nodos para asegurarse de que todos los sistemas físicos (*phys-schost-1*, *phys-schost-2*) actúen como maestros predeterminados de una aplicación. Cada sistema físico actúa como el nodo secundario del otro. Todos los proyectos asociados con las aplicaciones 1 y 2 deben estar representado en los archivos de la base de datos de proyectos de ambos nodos. Cuando el clúster se está ejecutando normalmente, todas las aplicaciones se ejecutan en su principal predeterminado, donde la función de gestión le asigna todo el tiempo de CPU.

Después de una recuperación de fallos o cambio, las dos aplicaciones se ejecutan en un único nodo en que se les asigna particiones, como se especifica en el archivo de configuración. Por ejemplo, esta entrada del archivo `/etc/project` especifica que a la aplicación 1 se le asignan 4 particiones y a la aplicación 2 se le asigna 1 partición.

```
Prj_1:100:project for App-1:root::project.cpu-shares=(privileged,4,none)
Prj_2:101:project for App-2:root::project.cpu-shares=(privileged,1,none)
```

El diagrama siguiente ilustra las operaciones normales y de recuperación de fallos de esta configuración. El número de particiones que se asignen no cambia. Sin embargo, el porcentaje de tiempo de CPU disponible para cada aplicación puede cambiar, dependiendo del número de particiones asignadas a cada proceso que requiera tiempo de CPU.



Clúster de dos nodos con tres aplicaciones

En un clúster de dos nodos con tres aplicaciones, se puede configurar un sistema físico (*phys-schost-1*) como principal predeterminado de una aplicación y el segundo sistema físico (*phys-schost-2*) como maestro predeterminado para las otras dos aplicaciones. Consideremos el siguiente archivo de base de datos de proyectos de ejemplo que está en todos los nodos. El archivo de base de datos de proyectos no cambia cuando se produce una recuperación de fallos o un cambio.

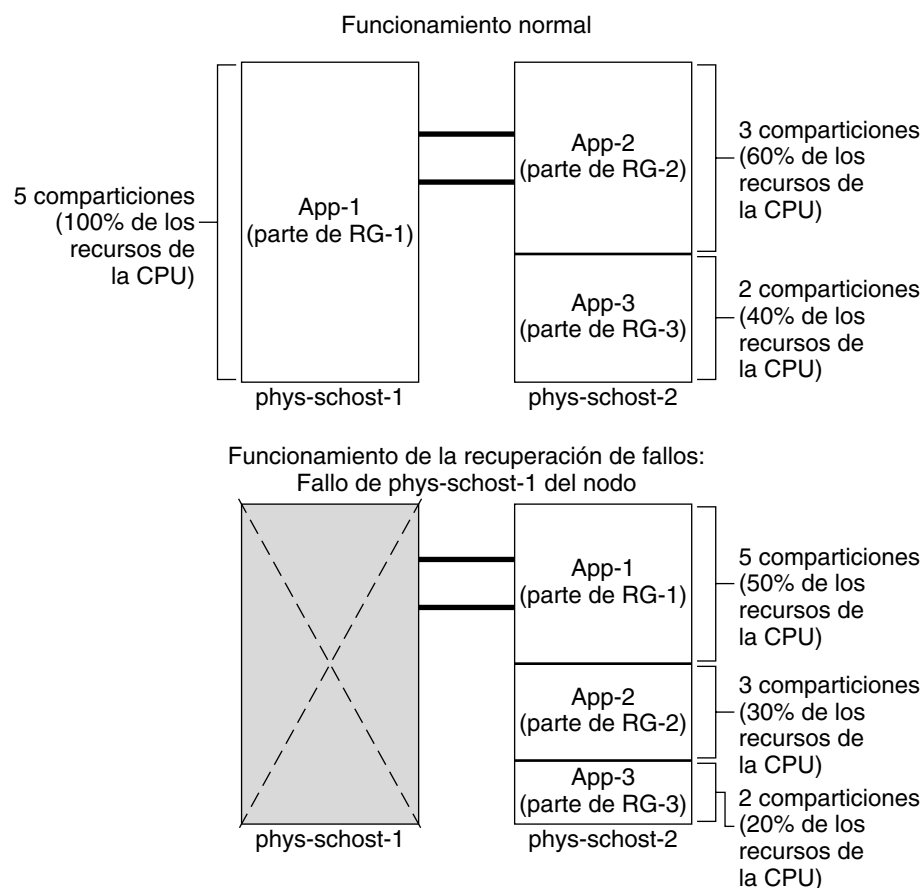
```
Prj_1:103:project for App-1:root::project.cpu-shares=(privileged,5,none)
Prj_2:104:project for App_2:root::project.cpu-shares=(privileged,3,none)
Prj_3:105:project for App_3:root::project.cpu-shares=(privileged,2,none)
```

Cuando el clúster se está ejecutando normalmente, a la aplicación 1 se le asignan 5 particiones en su maestro predeterminado, *phys-schost-1*. Este número es equivalente al 100 por cien del tiempo de CPU porque es la única aplicación que lo demanda en ese nodo. A las aplicaciones 2 y 3 se les asigna 3 y 2 particiones respectivamente en su maestro predeterminado, *phys-schost-2*. La aplicación 2 recibiría el 60 por ciento del tiempo de CPU y la aplicación 3 recibiría el 40 por ciento durante el funcionamiento normal.

Si se produce una recuperación de fallos o una conmutación y la aplicación 1 se cambia a *phys-schost-2*, las particiones para las tres aplicaciones siguen siendo las mismas. Sin embargo, los porcentajes de recursos de CPU se reasignan según el archivo de base de datos de proyectos.

- La aplicación 1, con 5 particiones, recibe el 50 por ciento de CPU.
- La aplicación 2, con tres particiones, recibe el 30 por ciento de CPU.
- La aplicación 3, con 2 particiones, recibe el 20 por ciento de CPU.

El diagrama siguiente ilustra el funcionamiento normal y las operaciones de recuperación de fallos de esta configuración.



Recuperación de fallos sólo de grupos de recursos

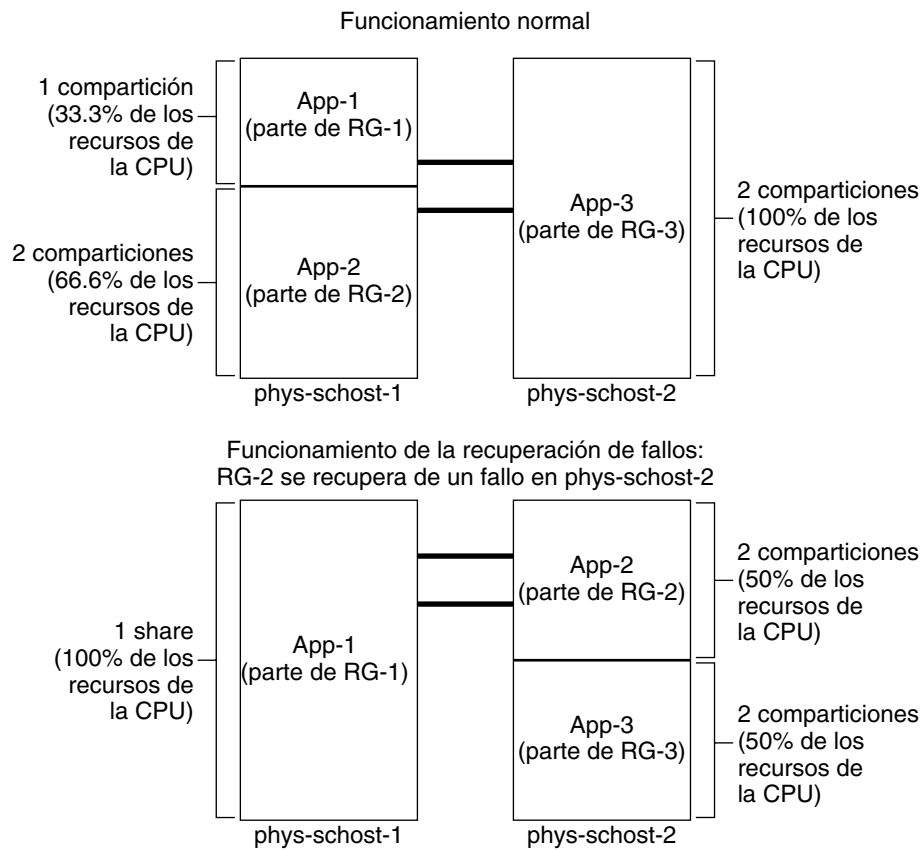
En una configuración en la que varios grupos de recursos tienen el mismo maestro predeterminado, un grupo de recursos (y sus aplicaciones asociadas) pueden entrar en recuperación de fallos o cambiarse a un nodo secundario. Mientras tanto el maestro predeterminado está ejecutándose en el clúster.

Nota – Durante la recuperación de fallos a la aplicación que falla se le asignan recursos, como los que se especifican en el archivo de configuración en el nodo secundario. En este ejemplo, los archivos de la base de datos de proyectos de los nodos primario y secundario tienen las mismas configuraciones.

Por ejemplo, este archivo de configuración de ejemplo especifica que a la aplicación 1 se le asigna 1 partición, a la aplicación 2 se le asignan 2 y a la aplicación 3 se le asignan otras 2.

```
Prj_1:106:project for App_1:root::project.cpu-shares=(privileged,1,none)
Prj_2:107:project for App_2:root::project.cpu-shares=(privileged,2,none)
Prj_3:108:project for App_3:root::project.cpu-shares=(privileged,2,none)
```

El diagrama siguiente ilustra las operaciones normal y de resolución de problemas de esta configuración, en que RG-2, que contiene la aplicación 2, falla sobre *phys-schost-2*. Tenga en cuenta que el número de particiones asignadas no cambia. Sin embargo, el porcentaje de tiempo de CPU disponible para cada aplicación puede cambiar, dependiendo del número de particiones asignadas a cada aplicación que requiera tiempo de CPU.



Adaptadores de red pública y IP Network Multipathing

Los clientes hacen peticiones de datos al clúster a través de la red pública. Cada nodo del clúster está conectado como mínimo a una red pública a través de un par de adaptadores.

El software de Solaris Internet Protocol (IP) Network Multipathing en Sun Cluster ofrece el mecanismo básico para supervisar los adaptadores de red pública y cambiar direcciones IP de un adaptador a otro cuando se detecta un fallo. Todos los nodos del clúster tienen su propia configuración IP Network Multipathing, que puede ser distinta de la de otros nodos del clúster.

Los adaptadores de red pública están organizados en *grupos de ruta múltiple IP* (grupos de ruta múltiple). Cada grupo de ruta múltiple tiene uno o más adaptadores de red pública. Cada adaptador de un grupo de ruta múltiple puede estar activo o se pueden configurar interfaces en espera que estén inactivos a menos que ocurra una recuperación de fallos. El daemon de ruta múltiple `in.mpathd` usa una dirección IP de prueba para detectar fallos y reparaciones, si detecta un fallo en uno de los adaptadores, se produce una recuperación de fallos. Todos los accesos de red se transfieren desde el adaptador erróneo a otro que funcione del grupo de ruta múltiple, con lo cual se mantiene la conectividad de red pública para el nodo. Si se configuró una interfaz en espera, el daemon lo elegirá. En caso contrario, `in.mpathd` elige la interfaz con la dirección IP de número inferior. Debido a que la recuperación del fallo se produce en la interfaz del adaptador, las conexiones de mayor nivel como TCP no se resultan afectadas excepto por un breve retraso durante la recuperación de fallos. Cuando se completa satisfactoriamente la recuperación de fallos de direcciones IP, se envían emisiones ARP generalizadas. Así se mantiene la conectividad con clientes remotos.

Nota – Debido a la característica de recuperación de congestiones de TCP, los puntos finales de TCP pueden verse más retrasados después de una recuperación de fallos satisfactoria, ya que algunos segmentos podrían perderse durante el proceso, activándose el mecanismo de control de congestión en TCP.

Los grupos de ruta múltiple ofrecen bloques de construcción para nombres de sistema lógicos y recursos de dirección compartida. También se pueden crear grupos de ruta múltiple independientemente de los nombres de sistema lógicos y los recursos de dirección compartida para supervisar la conectividad de red pública de los nodos del clúster. El mismo grupo de ruta múltiple de un nodo puede alojar cualquier número de nombres de sistema lógicos o recursos de dirección compartida. Si desea más información sobre sistemas lógicos y recursos de direcciones compartidos, consulte *Sun Cluster Data Services Planning and Administration Guide for Solaris OS*.

Nota – El diseño del mecanismo IP Network Multipathing está pensado para detectar y filtrar fallos de adaptadores, no para recuperarse de un administrador usando `ifconfig(1M)` para quitar una de las direcciones IP lógicas (o compartidas). El software de Sun Cluster trata las direcciones IP lógicas y compartidas como recursos gestionados por RGM. La forma correcta de que un administrador agregue o retire una dirección IP es usar `scrgadm(1M)` para modificar el grupo de recursos que lo contiene.

Para obtener más información sobre la implementación de Solaris de la Ruta múltiple de red IP, consulte la documentación apropiada para el sistema operativo Solaris instalado en su clúster.

Versión del sistema operativo	Si desea obtener más instrucciones, vaya a...
Sistema operativo Solaris 8	<i>IP Network Multipathing Administration Guide</i>
Sistema operativo Solaris 9	“IP Network Multipathing Topics” en <i>System Administration Guide: IP Services</i>

SPARC: Compatibilidad con la reconfiguración dinámica

La compatibilidad de Sun Cluster 3.1 4/04 con la función de software de reconfiguración dinámica (DR) se está desarrollando en fases incrementales. Este apartado describe los conceptos y consideraciones con respecto a la compatibilidad de Sun Cluster 3.1 4/04 con la función DR.

Tenga en cuenta que todos los requisitos, procedimientos y restricciones que se documentan para la función DR de Solaris[00c2][00a0] también se aplican al soporte DR de Sun Cluster (excepto para el funcionamiento silencioso del sistema operativo). Por consiguiente, se ha de repasar la documentación de la función de DR de Solaris *antes* de utilizarla con el software Sun Cluster. en concreto las cuestiones que afectan a los dispositivos de E/S que no son de la red durante una operación de desconexión de DR. Están disponibles para su descarga *Sun Enterprise 10000 Dynamic Reconfiguration User Guide* y *Sun Enterprise 10000 Dynamic Reconfiguration Reference Manual* (de las colecciones *Solaris 8 on Sun Hardware* o *Solaris 9 on Sun Hardware*) desde <http://docs.sun.com>.

SPARC: Descripción general de la reconfiguración dinámica

La función DR permite operaciones, como la extracción de hardware en sistemas que están en funcionamiento. Los procesos DR están diseñados para asegurar el funcionamiento continuo del sistema sin necesidad de pararlo o interrumpir la disponibilidad del clúster.

DR funciona a nivel de placa, por lo tanto, afecta a todos los componentes de ésta. Cada placa puede contener varios componentes, como CPU, memorias e interfaces periféricas para unidades de disco, unidades de cinta y conexiones en red.

Extraer una placa que contenga componentes activos puede provocar errores en el sistema. Antes de extraer una placa, el subsistema DR consulta otros subsistemas, como Sun Cluster, para determinar si los componentes de la placa se están utilizando. Si el subsistema DR encuentra que la placa se está usando, la operación de extraer la placa por DR no se lleva a cabo. Por tanto, siempre es seguro usar una operación de extracción de placa por DR ya que el subsistema DR rechaza las operaciones en placas que contengan componentes activos.

La operación de agregar placas DR también es siempre segura. El sistema pone en funcionamiento automáticamente las CPU y la memoria de una placa recién añadida. Sin embargo, el administrador del sistema debe configurar manualmente el clúster para usar activamente componentes que están en la placa recién añadida.

Nota – El subsistema DR tiene varios niveles. Si un nivel inferior informa de un error, el superior también informa del mismo error. Sin embargo, cuando el nivel inferior informa del error específico, el superior informará de un “error desconocido”. Los administradores del sistema deben obviar el “error desconocido” del que informa el nivel superior.

Los apartados siguientes describen consideraciones de DR para los distintos tipos de dispositivos.

SPARC: Consideraciones de la agrupación DR para dispositivos de CPU

El software Sun Cluster no rechazará una operación de extracción de placa con DR debido a la presencia de dispositivos de CPU.

Cuando una operación de añadir placa con DR tenga éxito, los dispositivos de CPU de la placa añadida se incorporarán automáticamente al funcionamiento del sistema.

SPARC: Consideraciones de la agrupación DR para memoria

Con respecto a DR, existen dos tipos de memoria que considerar que difieren sólo en su uso. El hardware real es el mismo para ambos tipos.

La memoria usada por el sistema operativo se llama caja de la memoria del núcleo. El software Sun Cluster no admite operaciones de extracción en placas que contengan la caja de la memoria del núcleo y rechazará cualquiera de estas operaciones. Cuando una operación de extracción de placa con DR pertenezca a una memoria distinta de la caja de la memoria del núcleo, Sun Cluster no rechazará la operación.

Cuando una operación de añadir placa con DR que pertenezca a memoria tenga éxito, la memoria de la placa añadida se incorporará automáticamente al funcionamiento del sistema.

SPARC: Consideraciones de la agrupación DR para unidades de disco y cinta

Sun Cluster rechaza las operaciones DR de extraer-placa en las unidades activas del nodo principal. Las operaciones de extracción de placa con DR pueden llevarse a cabo en unidades no activas del nodo primario y en cualquier unidad del nodo secundario. Después de la operación de DR, el acceso a los datos del clúster prosigue de la misma forma.

Nota – Sun Cluster rechaza las operaciones DR que afecten a la disponibilidad de los dispositivos del quórum. Para consideraciones sobre los dispositivos del quórum y el procedimiento para llevar a cabo operaciones DR en ellos, consulte “SPARC: Consideraciones de la agrupación DR para los dispositivos del quórum” en la página 92.

Consulte “ Task Map: Dynamic Reconfiguration with Quorum Devices” en *Sun Cluster System Administration Guide for Solaris OS* para obtener instrucciones detalladas sobre cómo realizar estas acciones.

SPARC: Consideraciones de la agrupación DR para los dispositivos del quórum

Si la operación de extracción de placa con DR pertenece a una placa que contenga una interfaz con un dispositivo configurado para el quórum, Sun Cluster rechazará la operación e identificará el dispositivo del quórum que podría resultar afectado por ésta. Es necesario inhabilitar el dispositivo como del quórum antes de poder efectuar una operación de extracción de placa con DR.

Consulte “ Task Map: Dynamic Reconfiguration with Quorum Devices” en *Sun Cluster System Administration Guide for Solaris OS* para obtener instrucciones detalladas sobre cómo realizar estas acciones.

SPARC: Consideraciones de la agrupación DR para interfaces de interconexión del clúster

Si la operación de extracción de placa con DR pertenece a una placa que contenga una interfaz de interconexión del clúster activa, Sun Cluster rechazará la operación e identificará la interfaz que podría resultar afectada por la operación. Es necesario usar una herramienta de administración de Sun Cluster para inhabilitar la interfaz activa antes de que la operación de DR pueda tener éxito (consulte también la nota de precaución siguiente).

Consulte “ Administering the Cluster Interconnects” in *Sun Cluster System Administration Guide for Solaris OS* para obtener instrucciones detalladas sobre cómo realizar estas acciones.



Caution – Sun Cluster requiere que todos los nodos del clúster tengan al menos una ruta de acceso que funcione con cada uno de los demás nodos del clúster. No inhabilite una interfaz de interconexión privada en el caso de que represente la última ruta a cualquiera de los nodos del clúster.

SPARC: Consideraciones de la agrupación DR para interfaces de red pública

Si la operación de extracción de placa con DR pertenece a una placa que contenga una interfaz de red pública activa, Sun Cluster rechazará la operación e identificará la interfaz que podría resultar afectada por la operación. Antes de suprimir una placa que contenga una interfaz de red activa, todo el tráfico de esa interfaz debe cambiarse a otra interfaz funcional del grupo de ruta múltiple mediante la orden `if_mpadm(1M)`.



Caution – Si el resto de adaptadores de red fallan mientras se está efectuando una supresión de DR en el adaptador de red inhabilitado, la disponibilidad se verá afectada. El adaptador restante no tiene a quién transferir el control durante la operación de DR.

Consulte “ Administering the Public Network” en *Sun Cluster System Administration Guide for Solaris OS* para obtener instrucciones detalladas sobre cómo realizar una operación de eliminación de DR en una interfaz de red pública.

Preguntas más frecuentes (FAQ)

INDEXTERM-343

Este capítulo incluye las respuestas a las preguntas más frecuentes sobre el sistema SunPlex. Las preguntas están organizadas por temas.

FAQ sobre alta disponibilidad

- **¿Qué es exactamente un sistema de alta disponibilidad?**

El sistema SunPlex define la alta disponibilidad (HA) como la posibilidad del clúster de mantener en marcha una aplicación aunque se produzca un fallo que en condiciones normales provocaría que el servidor no estuviera disponible.

- **¿Cuál es el proceso por el que el clúster proporciona alta disponibilidad?**

A través de un proceso conocido como recuperación de fallos, la estructura del clúster proporciona un entorno de alta disponibilidad. Recuperación de fallos es una serie de pasos que realiza el clúster para migrar recursos de servicios de datos de un nodo fallido a otro operativo dentro del clúster.

- **¿Cuál es la diferencia entre un servicio de datos a prueba de fallos y otro escalable?**

Existen dos tipos de servicios de datos de alta disponibilidad, a prueba de fallos y escalable.

Un servicio de datos a prueba de fallos ejecuta una aplicación sólo en un nodo primario del clúster cada vez. Los demás nodos pueden ejecutar otras aplicaciones, pero cada una se ejecuta sólo en un nodo. Si un nodo primario falla, las aplicaciones que se ejecutan en este nodo se trasladan a otro nodo y continúan ejecutándose.

Un servicio escalable reparte una aplicación entre varios nodos para crear un servicio lógico único que aprovecha el número de nodos y procesadores de todo el clúster en el que se ejecutan.

Para cada aplicación un nodo aloja la interfaz física para el clúster. Este nodo se denomina interfaz global (GIF). Pueden haber varios nodos GIF en el clúster. Cada uno de ellos aloja una o más interfaces lógicas que pueden usar los servicios escalables y que reciben el nombre de *interfaces globales*. Un nodo GIF aloja una interfaz global para todas las solicitudes hacia una aplicación en particular y las despacha a los distintos nodos en los que se esté ejecutando el servidor de la aplicación. Si el nodo GIF falla, la interfaz global se traslada a un nodo superviviente.

Si falla alguno de los nodos en los que se está ejecutando la aplicación, ésta continúa ejecutándose en los demás nodos con alguna pérdida de rendimiento hasta que el nodo fallido vuelve al clúster.

FAQ sobre sistemas de archivos

- **¿Puedo ejecutar uno o más nodos del clúster como servidores NFS de alta disponibilidad con otros nodos de clúster como clientes?**

No, no haga un montaje de bucle cerrado.

- **¿Puedo usar un sistema de archivos del clúster para aplicaciones que no estén bajo el control del Resource Group Manager?**

Sí. Sin embargo, sin el control del RGM es necesario reiniciar las aplicaciones manualmente después que falle el nodo en el que se están ejecutando.

- **¿Deben tener todos los sistemas de archivos del clúster un punto de montaje debajo del directorio /global?**

No. Sin embargo, si se sitúan los sistemas de archivos clúster bajo el mismo punto de montaje, como /global, se consigue una mejor organización y gestión de estos sistemas de archivos.

- **¿Cuáles son las diferencias entre usar el sistema de archivos del clúster y exportar sistemas de archivos NFS?**

Existen varias diferencias:

1. El sistema de archivos del clúster admite dispositivos globales. NFS no admite acceso remoto a los dispositivos.
2. El sistema de archivos del clúster dispone de un espacio de nombres global. Sólo es necesaria una orden de montaje. Con NFS, debe montarse el sistema de archivos en cada uno de los nodos.
3. El sistema de archivos del clúster almacena en antememoria los archivos en más casos que NFS. Por ejemplo, cuando se accede a un archivo desde varios nodos, para lectura, escritura, bloqueo de archivo, E/S asíncrona.

4. El sistema de archivos del clúster está creado para aprovechar futuras interconexiones rápidas de clúster que ofrezcan DMA remotas y funciones de copia cero.
 5. Si modifica los atributos de un archivo (por ejemplo, mediante `chmod (1M)`) en un sistema de archivos del clúster, los cambios se reflejarán inmediatamente en todos los nodos. Con un sistema de archivos NFS exportado, esto puede llevar mucho más tiempo.
- **El sistema de archivos `/global/.devices/node@<ID_nodo>` aparece en los nodos de mi clúster. ¿Puedo usar este sistema de archivos para almacenar datos que deseo que estén altamente disponibles y sean globales?**

Estos sistemas de archivos almacenan los espacios de nombre de dispositivo global. No están pensados para un uso general. Aunque sean globales, nunca se accede a ellos de forma global (cada nodo sólo accede a su propio espacio de nombres de dispositivo global). Si un nodo está caído, el resto no puede acceder al espacio de nombres del nodo en cuestión. Estos sistemas de archivos no son de alta disponibilidad. No deberían usarse para almacenar datos que hayan de estar globalmente accesibles o altamente disponibles.

FAQ de gestión de volúmenes

- **¿Necesito duplicar todos los dispositivos de disco?**

Para que un dispositivo de disco se considere de alta disponibilidad, debe estar duplicado o usar hardware RAID-5. Todos los servicios de datos deben usar dispositivos de disco de alta disponibilidad o sistemas de archivos del clúster montados en dispositivos de disco de alta disponibilidad. Estas configuraciones pueden tolerar fallos de disco puntuales.
- **¿Puedo usar un gestor de volúmenes para los discos locales (disco de arranque) y otro distinto para los discos multisistema?**

SPARC: Esta configuración se admite con el software Solaris Volume Manager cuando gestiona los discos locales y VERITAS Volume Manager con los discos multisistema. No se admite ninguna otra combinación.

x86: No, esta configuración no es compatible ya que sólo se admite Solaris Volume Manager en clústeres basados en x86.

FAQ de servicios de datos

- **¿Qué servicios de datos SunPlex están disponibles?**

La lista de servicios de datos admitidos se incluye en “Supported Products” en *Sun Cluster 3.1 9/04 Release Notes for Solaris OS* .

- **¿Qué versiones de aplicaciones admiten los servicios de datos de SunPlex?**

La lista de versiones de aplicaciones admitidas se incluye en “Supported Products” en *Sun Cluster 3.1 9/04 Release Notes for Solaris OS*.

- **¿Puedo escribir mi propio servicio de datos?**

Sí. Consulte “Data Service Development Library Reference” en *Sun Cluster Data Services Developer’s Guide for Solaris OS* para obtener más información.

- **Al crear recursos de red, ¿debo especificar direcciones IP numéricas o nombres de sistema?**

El método preferido para especificar recursos de red es usar el nombre de sistema de UNIX en lugar de la dirección IP numérica.

- **Al crear recursos de red, ¿cuál es la diferencia entre usar un nombre de sistema lógico (un recurso LogicalHostname) o una dirección compartida (un recurso SharedAddress)?**

Excepto en el caso de Sun Cluster HA for NFS, siempre que en la documentación se habla del uso de un recurso LogicalHostname de un grupo de recursos de modo Failover, en su lugar se puede usar un recurso SharedAddress o LogicalHostname. El uso de un recurso SharedAddress produce una sobrecarga adicional porque el software de red del clúster está configurado para SharedAddress pero no para LogicalHostname.

Existe una ventaja respecto a usar SharedAddress cuando que se está configurando servicios de datos escalables y a prueba de fallos y se desea que los clientes puedan acceder a ambos servicios usando el mismo nombre de sistema. En este caso, los recursos SharedAddress junto con el recurso de la aplicación de recuperación de fallos están contenidos en un grupo de recursos, mientras que el recurso del servicio escalable está contenido en un grupo de recursos separado y configurado para usar SharedAddress. Tanto los servicios escalables como los a prueba de fallos pueden usar después el mismo conjunto de nombres de sistema/direcciones que estén configurados en el recurso SharedAddress.

FAQ de redes públicas

- **¿Qué adaptadores de red pública admite el sistema SunPlex?**

Actualmente el sistema SunPlex admite adaptadores de red pública Ethernet (10/100BASE-T y 1000BASE-SX Gb). Debido a que en el futuro podrían admitirse interfaces nuevas, compruebe con el representante de ventas de Sun cuál es la información más actual.

- **¿Cuál es el papel de la dirección MAC en la recuperación de fallos?**

Cuando se produce una recuperación de fallos, se generan nuevos paquetes de protocolo de resolución de direcciones (ARP) y se envían a todo el mundo. Estos paquetes ARP contienen la nueva dirección MAC (del nuevo adaptador físico al que se ha transferido el nodo) y la dirección IP antigua. Cuando otra máquina de la red recibe uno de estos paquetes, borra la correlación MAC-IP antigua de la antememoria ARP y usa la nueva.

- **¿El sistema SunPlex admite el valor `local-mac-address?=true`?**

Sí. Además, la Ruta múltiple de red IP requiere que el valor `local-mac-address?` esté configurado como `true`.

Se puede establecer `local-mac-address?` con `eeprom(1M)`, en el indicador `ok` de la PROM de OpenBoot en un clúster basado en la plataforma SPARC o con la utilidad SCSI que se puede ejecutar opcionalmente tras un arranque de la BIOS en un clúster basado en la plataforma x86.

- **¿Cuánto retraso puede esperarse cuando la IP Network Multipathing lleva a cabo una conmutación entre adaptadores?**

El retraso podría ser de varios minutos. Esto es debido a que cuando se lleva a cabo el cambio de IP Network Multipathing, implica enviar una ARP innecesaria. Sin embargo, no hay garantías de que el encaminador entre el cliente y el clúster utilizará la ARP innecesaria. Así que, hasta que la entrada de la antememoria ARP de esta dirección IP en el encaminador no supere el tiempo de espera, es posible que utilice la dirección MAC anterior.

- **¿Con qué velocidad se detecta el fallo de un adaptador de red?**

El tiempo de detección de fallos predeterminado es de 10 segundos. Este algoritmo intenta cumplir el tiempo de detección de fallos predeterminado, pero el real depende de la carga de la red.

FAQ sobre pertenencia al clúster

- **¿Deben de tener todos los miembros del clúster la misma contraseña de superusuario?**

No es necesario que el superusuario comparta la misma contraseña en todos los miembros del clúster. Sin embargo, esta práctica simplifica mucho la administración del clúster.

- **¿Es importante el orden en el que arrancan los nodos?**

En la mayoría de los casos, no. Sin embargo es importante para evitar la amnesia (consulte “[Aislamiento de fallos](#)” en la [página 55](#) para obtener más información sobre la amnesia). Por ejemplo, si el nodo dos era el propietario del dispositivo del quórum y el nodo uno estaba caído, y después se desactiva el nodo dos, hay que volver a activar éste antes que aquél. Esto evita que accidentalmente se active un nodo con información de configuración del clúster anticuada.

- **¿Es necesario realizar una duplicación de los discos locales en un nodo del clúster?**

Sí. Aunque esta duplicación no es un requisito, con ella se asegura que un fallo de disco sin duplicación haga inoperativo el nodo. El inconveniente de realizar una duplicación de los discos locales de un nodo del clúster es que complica la administración del sistema.

- **¿Cuáles son los inconvenientes de realizar copias de respaldo de los miembros del clúster?**

En un clúster se pueden usar varios métodos de copia de respaldo. Un método es tener un nodo como respaldo con una unidad de cinta/biblioteca conectada. A continuación se debe usar el sistema de archivos del clúster para realizar la copia de seguridad de los datos. No conecte este nodo a los discos compartidos.

Consulte “[Backing Up and Restoring a Cluster](#)” en *Sun Cluster System Administration Guide for Solaris OS* para obtener información adicional sobre cómo realizar copias de seguridad y restaurar los datos.

- **¿Cuándo está un nodo en el suficiente buen estado como para usarlo como secundario?**

Después de un arranque, un nodo está en el suficiente buen estado como para ser secundario cuando éste muestra el indicador de inicio de sesión.

FAQ de almacenamiento del clúster

- **¿Qué hace que el almacenamiento multisistema sea de alta disponibilidad?**

El almacenamiento multisistema es de alta disponibilidad porque puede sobrevivir a la pérdida de un disco, gracias a la duplicación (o debido a controladores RAID-5 basados en el hardware). Como los dispositivos de almacenamiento multisistema tienen más de una conexión de sistema, también pueden resistir la pérdida de un nodo al cual estén conectados. Además, las rutas redundantes desde cada nodo al almacenamiento conectado ofrecen tolerancia para el fallo de un adaptador del bus del sistema, de un cable o del controlador de disco.

FAQ de interconexión del clúster

- **¿Qué interconexiones del clúster admite el sistema SunPlex?**

Actualmente, el sistema SunPlex admite las interconexiones del clúster Ethernet (100BASE-T Fast Ethernet y 1000BASE-SX Gb) en los clústers basados en las plataformas SPARC y x86. El sistema SunPlex admite las interconexiones del clúster de interfaz de red SCI sólo en los clústers basados en la plataforma SPARC.

- **¿Cuál es la diferencia entre un “cable” y una “ruta” de transporte?**

Los cables de transporte del clúster se configuran mediante los adaptadores de transporte y los conmutadores. Los cables unen adaptadores y conmutadores según cada componente. El gestor de la topología del clúster usa los cables disponibles para crear las rutas de transporte entre los extremos de los nodos. Un cable no se correlaciona directamente con una ruta de transporte.

Los cables son “habilitados” e “inhabilitados” por el administrador estáticamente. Los cables tienen un “estado,” (activados o desactivados) pero no un “estado.” Si se inhabilita un cable, es como si estuviera sin configurar. Los cables que están inhabilitados no pueden usarse como rutas de transporte. No han sido sondeados y por tanto no es posible conocer su estatus. El estado de un cable puede verse con `scconf - p`.

El gestor de topología del clúster establece dinámicamente las rutas de transporte. El “estatus” de una ruta de transporte viene determinado por el gestor de topologías. Una ruta puede tener un estado de “en línea” o “fuera de línea.” El estatus de una ruta de transporte puede visualizarse con `scstat (1M)`.

Considere el ejemplo siguiente de clúster de dos nodos con cuatro cables.

```
node1:adapter0      to switch1, port0
node1:adapter1      to switch2, port0
node2:adapter0      to switch1, port1
```

```
node2:adapter1    to switch2, port1
```

A partir de estos cuatro cables se pueden formar dos posibles rutas de transporte.

```
node1:adapter0    to node2:adapter0
```

```
node2:adapter1    to node2:adapter1
```

FAQ sobre sistemas cliente

- **¿Es necesario considerar alguna necesidad o restricción de cliente para usarla con un clúster?**

Los sistemas cliente se conectan al clúster como lo harían con cualquier otro servidor. En algunos casos, dependiendo de la aplicación del servicio de datos, es posible que necesite instalar software de lado del cliente o llevar a cabo otros cambios de configuración para que el cliente pueda conectarse a la aplicación del servicio de datos. Consulte los capítulos de *Sun Cluster Data Services Planning and Administration Guide* relacionados con este tema para obtener más información sobre los requisitos de la configuración del lado del cliente.

FAQ de consola de administración

- **¿Requiere el sistema SunPlex una consola de administración?**

Sí.

- **¿Tiene que tratarse de una consola de administración exclusiva para el clúster? ¿O puede usarse para otras tareas?**

El sistema SunPlex no requiere el uso de una consola de administración exclusiva, pero si se utiliza una de este tipo se obtienen estas ventajas:

- Permite la gestión centralizada del clúster ya que agrupa herramientas de consola y gestión en la misma máquina
- Ofrece al proveedor de servicio de hardware una resolución de problemas potencialmente mas rápida

- **¿Es necesario que la consola de administración se encuentre “cerca” del propio clúster, por ejemplo en la misma habitación?**

Compruébelo con su proveedor de servicio de hardware. Es posible que el proveedor le requiera que la consola se sitúe muy cerca del clúster en sí mismo. No existe ninguna razón técnica por la que la consola deba estar situada en la misma habitación.

- **¿Puede servirse a más de un clúster desde la misma consola de administración, siempre que se cumplan los requisitos de distancia?**

Sí. Se pueden controlar varios clústers desde una misma consola de administración. También puede compartirse un concentrador de terminal simple entre varios clústers.

FAQ sobre concentrador de terminal y procesador de servicio del sistema

- **¿Requiere el sistema SunPlex un concentrador de terminal?**

Todas las versiones del software empezando por Sun Cluster 3.0 no requieren la ejecución de ningún concentrador de terminal. Al contrario de lo que ocurre con el producto Sun Cluster 2.2, que requería un concentrador de terminal para el aislamiento de fallos, los productos posteriores no dependen del concentrador de terminal.

- **Veo que la mayoría de servidores de SunPlex usan concentradores de terminal, pero el Sun Enterprise E10000 server no. ¿Por qué?**

El concentrador de terminal es en la práctica un conversor serie para Ethernet para la mayoría de los servidores. Su puerto de consola es un puerto serie. El Sun Enterprise E10000 server no dispone de consola serie. El procesador de servicio del sistema (SSP) es la consola, a través de Ethernet o del puerto jtag. Para el Sun Enterprise E10000 server, siempre se usa SSP para consolas.

- **¿Cuáles son las ventajas de usar un concentrador de terminal?**

Usar un concentrador de terminal ofrece acceso en el nivel de la consola a todos los nodos desde una estación de trabajo remota de la red, incluso cuando el nodo está en la PROM de OpenBoot (OBP) en un nodo basado en la plataforma SPARC o en un subsistema de arranque en un nodo basado en la plataforma x86.

- **Si utilizo un concentrador de terminal no admitido por Sun, ¿qué necesito saber para certificar el que deseo utilizar?**

La diferencia principal entre el concentrador de terminal que admite Sun y los otros dispositivos de consola es que el de Sun dispone de firmware especial que evita que éste envíe un carácter de interrupción a la consola cuando arranca. Tenga en cuenta que si dispone de un dispositivo de consola que envía a la consola caracteres de interrupción o una señal que pueda interpretarse como tales, el nodo se apaga.

- **¿Se puede liberar un puerto bloqueado en el concentrador de terminal que admite Sun sin volverlo a arrancar?**

Sí. Anote el número de puerto que necesita restaurarse y escriba las órdenes siguientes:

```
telnet ct
Enter Annex port name or number: cli
annex: su -
```

```
annex# admin
admin : reset número_puerto
admin : quit
annex# hangup
#
```

Consulte los siguientes manuales para obtener más información sobre cómo configurar y administrar el concentrador de terminales compatible con Sun.

- “Administering Sun Cluster Overview” en *Sun Cluster System Administration Guide for Solaris OS*
- “Installing and Configuring the Terminal Concentrator” en *Sun Cluster 3.x Hardware Administration Manual for Solaris OS*
- **¿Qué ocurre si falla el propio concentrador de terminal? ¿Es necesario tener otro de recambio?**

No. No se pierde ninguna disponibilidad del clúster si el concentrador de terminal falla. Se pierde la posibilidad de conectarse a las consolas del nodo hasta que el concentrador vuelva a estar operativo.

- **Si se usa un concentrador de terminal, ¿qué ocurre con la seguridad?**

En general, el concentrador de terminal está conectado a una pequeña red que usan los administradores de sistema, no una red que se use para otros accesos de clientes. La seguridad se puede controlar limitando el acceso a esa red en particular.
- **SPARC: ¿Cómo se utiliza la reconfiguración dinámica con una cinta o unidad de disco?**
 - Determine si el disco o la unidad de cinta forma parte de un grupo de dispositivos activo. Si la unidad no forma parte de un grupo de dispositivos activos, puede llevar a cabo la operación de extracción DR en él.
 - Si la operación de extracción de placa DR pudiera afectar a un disco o unidad de cinta activos, el sistema rechazará la operación e identificará las unidades que se verían afectadas por la operación. Si la unidad forma parte de un grupo de dispositivos activos, vaya a [“SPARC: Consideraciones de la agrupación DR para unidades de disco y cinta” en la página 92.](#)
 - Determine si la unidad es un componente del nodo primario o del secundario. Si la unidad es un componente del nodo secundario, puede llevar a cabo la operación de extracción DR en el mismo.
 - Si la unidad es un componente del nodo primario, es necesario intercambiar los nodos primario y secundario antes de realizar la operación de extracción DR en el dispositivo.



Caution – Si el nodo principal falla durante una operación de DR en un nodo secundario, la disponibilidad del clúster queda afectada. El nodo primario no tiene dónde transferirse hasta que se proporcione un nodo secundario nuevo.

Índice

A

- a prueba de fallos
 - FAQ, 95
 - frente a escalable, 95
 - servicios de datos, 68
- adaptadores, *Ver* red, adaptadores
- Administración, clúster, 37-93
- agentes, *Ver* servicios de datos
- Aislamiento, 40
- aislar, 55
- almacenamiento, 24
 - SCSI, 25
- almacenar
 - FAQ, 101
 - reconfiguración dinámica, 92
- alta disponibilidad
 - Ver también* altamente disponible
 - estructura, 39
 - FAQ, 95
- altamente disponible
 - Ver también* alta disponibilidad
 - servicios de datos, 40
- Amnesia, 54
- apagado, 40
- API, 72, 77
- API DSDL, 77
- aplicación, *Ver* servicios de datos
- Aplicaciones de misión crítica, 63
- atributos, *Ver* propiedades
- auto-boot?, parámetro, 40
- aviso grave, 40, 41, 57

B

- bloqueo de archivo, 47

C

- cable, transporte, 101
- CCP, 29
- CCR, 41
- Clúster
 - administración, 37-93
- clúster
 - administrar, 37
 - componentes del software, 23
 - configurar, 41
 - Solaris Resource Manager, 79
 - contraseña, 100
- Clúster
 - desarrollo de aplicaciones, 37-93
- clúster
 - descripción, 14
 - extracción de placa, 92
 - FAQ de almacenamiento, 101
 - hardware, 15, 21
 - hora, 38
 - interconexión, 22, 27
 - adaptadores, 27
 - admitida, 101
 - cables, 27
 - FAQ, 101
 - interfaces, 27
 - reconfiguración dinámica, 93
 - servicios de datos, 74

- clúster, interconexión (Continuación)
 - uniones, 27
 - interfaz de red pública, 66
 - lista de tareas, 19
 - medios, 26
 - miembros, 22, 40
 - FAQ, 100
 - reconfiguración, 40
 - nodos, 22
 - objetivos, 14
 - orden de arranque, 100
 - punto de vista del administrador de sistemas, 16
 - punto de vista del programador de aplicaciones, 17
 - red pública, 27
 - reparación, 15
 - respaldo, 100
 - servicios de datos, 65
 - sistema de archivos, 47
 - HASStoragePlus, 49
 - uso, 48
 - sistemas de archivos, 96
 - FAQ
 - Ver también* sistemas de archivos
 - topologías, 29, 34
 - ventajas, 14
- Clústeres de aplicaciones reales de Oracle, 73
- CMM, 40
 - mecanismo de recuperación rápida, 40
 - Ver también* recuperación rápida
- componentes del software, 23
- concentrador de terminal, FAQ, 103
- configuración, límites de memoria
 - virtual, 82-83
- configuración cliente/servidor, 65
- Configuraciones, quórum, 58
- configuraciones de bases de datos paralelas, 22
- configurar
 - base datos paralela, 22
 - cliente/servidor, 65
 - depósito, 41
 - servicios de datos, 79
- conflicto de reserva, 56
- consola
 - acceso, 28
 - administración, 29
 - FAQ, 102

- consola (Continuación)
 - administrativa, 28
 - procesador de servicio del sistema, 28
- consola de administración, 29
 - FAQ, 102
- contraseña, root, 100
- contraseña de root, 100
- controlador, ID de dispositivo, 42

D

- datos, almacenar, 96
- depósito de configuración del clúster, 41
- desarrollo de aplicaciones, 72
- Desarrollo de aplicaciones, 37-93
 - /dev/global/espacio de nombres, 46
- DID, 42
- dirección compartida, 66
 - frente a nombre de sistema lógico, 98
 - nodo de interfaz global, 66
 - servicios de datos escalables, 68
- dirección IP, 98
- dirección MAC, 99
- disco de arranque, *Ver* discos, local
- discos
 - aislamiento de fallos, 55
 - dispositivos globales, 41, 46
 - dispositivos SCSI, 25
 - grupos de dispositivos, 43
 - multipuerto, 44
 - propiedad principal, 44-46
 - recuperación de fallos, 43
 - local, 26, 41, 46
 - duplicación, 100
 - gestión de volúmenes, 97
 - multisistema, 41, 43, 46
 - reconfiguración dinámica, 92
- discos locales, 26
- dispositivo
 - global, 41
 - ID, 42
 - multisistema, 24
- dispositivo multisistema, *Ver* dispositivos, multisistema
- Dispositivos, quórum, 53
- Distribución de aplicaciones, 59
- DR, *Ver* reconfiguración dinámica

E

E10000, *Ver* Sun Enterprise E10000

equilibrio de cargas, 70

escalable

FAQ, 95

frente a prueba de fallos, 95

escalables

grupos de recursos, 68

servicios de datos, 68

espacio de nombres

correlación, 47

global, 46

local, 47

Esquizofrenia, 54

esquizofrenia, aislamiento de fallos, 55

estructura, alta disponibilidad, 39

extracción de placas, reconfiguración

dinámica, 92

F

Fallo, aislamiento, 40

fallo

aislar, 55

detección, 39

recuperación, 39

retroceso, 71

FAQ, 95

a prueba de fallos frente a escalable, 95

almacenamiento del clúster, 101

alta disponibilidad, 95

concentrador de terminal, 103

consola de administración, 102

gestión de volúmenes, 97

interconexión del clúster, 101

miembros del clúster, 100

procesador de servicio de sistema, 103

red pública, 99

servicios de datos, 98

sistemas cliente, 102

sistemas de archivos, 96

G

gestión de volúmenes, dispositivos

multisistema, 25

gestión de recursos, 79

gestión de volúmenes

discos locales, 97

discos multisistema, 97

espacio de nombres, 46

FAQ, 97

gestor de volúmenes VERITAS, 97

RAID-5, 97

Solaris Volume Manager, 97

gestor de grupos de recursos, *Ver* RGM

global

dispositivo, 41, 43

discos locales, 26

montaje, 47

espacio de nombres, 42, 46

discos locales, 26

interfaz, 66, 95

servicios escalables, 68

/global, punto de montaje, 96

/global punto de montaje, 47

grupo de dispositivos, 43

cambiar propiedades, 44-46

grupos

dispositivo de discos

Ver discos, grupos de dispositivos

grupos de dispositivos de disco

multipuerto, 44

grupos de recursos, 76

a prueba de fallos, 68

configurar, 77

estados, 77

propiedades, 79

H

HA, *Ver* alta disponibilidad

hardware, 15, 21, 90

Ver también almacenamiento

Ver también discos

componentes de interconexión del

clúster, 27

fallo, 39

reconfiguración dinámica, 90

recuperación, 39

HASStoragePlus, 76

tipo de recurso, 49

hora, entre nodos, 38

I

ID

- dispositivo, 42
- nodo, 46

iniciador múltiple SCSI, 25

interfaces

- Ver red, interfaces*
- administrativas, 37

interfaces administrativas, 37

ioctl, 56

IPMP, *Ver Ruta múltiple de red IP*

L

local_mac_address, 99

LogicalHostname, *Ver nombre de sistema lógico*

M

medios, 26

- extraíbles, 26

modelo de servidor en clúster, 65

modelo de servidor individual, 65

montaje

- con syncdir, 49
- dispositivos globales, 47
- /global, 96
- sistemas de archivos, 47

N

Network Time Protocol, 38

NFS, 49

nodo de interfaz global, *Ver nodo de interfaz global*

nodo de respaldo, 100

nodo GIF, 95

nodo primario, 66

nodo secundario, 66

nodos, 22

- IDnodo, 46
- interfaz global, 66
- orden de arranque, 100
- primario, 66
- principales, 44-46

nodos (Continuación)

respaldo, 100

secundario, 66

secundarios, 44-46

nombre de sistema, 65

nombre de sistema lógico, 66

frente a dirección compartida, 98

servicios de datos a prueba de fallos, 68

NTP, 38

O

Oracle Parallel Server, *Ver Oracle Real*

Application Clusters

orden de arranque, 100

P

panel de control del clúster, 29

pertenencia, *Ver clúster, miembros*

Preguntas más frecuentes, *Ver FAQ*

procesador de servicio del sistema, 28, 29

procesador de sistema de servicio, FAQ, 103

programador, aplicaciones para clúster, 17

propiedad principal, grupos de dispositivos de discos, 44-46

Propiedad scsi-initiator-id, 26

propiedades

cambiar, 44-46

grupos de recursos, 79

recursos, 79

Resource_project_name, 81-82

RG_project_name, 81-82

proyectos, 79

proyectos de Solaris, 79

Q

Quórum, 53

configuraciones, 57, 58

configuraciones atípicas, 63

configuraciones recomendadas, 60-63

dispositivos, 53

- quórum
 - dispositivos
 - reconfiguración dinámica, 92
- Quórum
 - malas configuraciones, 64-65
- quórum
 - mejores prácticas, 58
 - recuentos de votos, 55
- Quórum
 - requisitos, 58

R

- reconfiguración dinámica, 90
 - descripción, 91
 - discos, 92
 - dispositivos de CPU, 91
 - dispositivos del quórum, 92
 - interconexión del clúster, 93
 - memoria, 92
 - red pública, 93
 - unidades de cinta, 92
- Recuentos de votos
 - dispositivos de quórum, 55
 - nodos, 55
- recuperación, 39
 - retroceso, 71
- recuperación de fallos
 - casos
 - Solaris Resource Manager, 83-89
 - grupos de dispositivos de disco, 43
- recuperación rápida, 40
 - aislamiento de fallos, 56
- recursos, 76
 - estados, 77
 - propiedades, 79
- red
 - adaptadores, 27, 89-90
 - dirección compartida, 66
 - equilibrio de cargas, 70
 - interfaces, 27, 89-90
 - nombre de sistema lógico, 66
 - privada
 - Ver clúster, interconexión
 - pública, 27
 - FAQ, 99
 - interfaces, 99

- red, pública (Continuación)
 - reconfiguración dinámica, 93
 - Ruta múltiple de red IP, 89-90
 - recursos, 66, 76
- red privada, Ver clúster, interconexión
- red pública, Ver red, pública
- reserva de grupo persistente, 56
- Resource_project_name, propiedad, 81-82
- resources, configurar, 77
- respaldo, 100
- retroceso, 71
- RG_project_name, propiedad, 81-82
- RGM, 67, 76, 79
- RMAPI, 77
- ruta de acceso, transporte, 101
- ruta múltiple, 89-90
- Ruta múltiple de red IP, 89-90
 - tiempo de recuperación de fallos, 99

S

- SCSI
 - aislamiento de fallos, 55
 - conflicto de reserva, 56
 - iniciador múltiple, 25
 - reserva de grupo persistente, 56
- servicio de datos, a prueba de fallos, 68
- servicios de datos, 65, 66
 - admitidos, 98
 - altamente disponible, 40
 - API, 72
 - API de biblioteca, 73
 - configurar, 79
 - desarrollar, 72
 - escalables, 68
 - FAQ, 98
 - grupos de recursos, 76
 - interconexión del clúster, 74
 - métodos, 67
 - recursos, 76
 - supervisor de fallos, 72
 - tipos de recursos, 76
- servidor
 - modelo de servidor en clúster, 65
 - modelo de servidor individual, 65
- SharedAddress, Ver dirección compartida

- sistema de archivos
 - clúster, 47
 - local, 49
 - montaje, 47
 - NFS, 49
 - syncdir, 49
 - UFS, 49
 - VxFS, 49
- sistema de archivos local, 49
- sistemas cliente, 28
 - restricciones, 102
- sistemas clientes, FAQ, 102
- sistemas de archivos
 - almacenar datos, 96
 - alta disponibilidad, 96
 - clúster, 96
 - FAQ, 96
 - global, 96
 - montaje, 96
 - NFS, 96
 - sistemas de archivos del clúster, 96
 - uso, 48
- software
 - fallo, 39
 - recuperación, 39
- Solaris Resource Manager, 79
 - casos de recuperación de fallos, 83-89
 - configuración de los límites de la memoria virtual, 82-83
 - requisitos de configuración, 81-82
- Solaris Volume Manager, dispositivos multisistema, 25
- SSP, *Ver* procesador de servicio del sistema
- Sun Cluster
 - Ver* clúster
- Sun Enterprise E1000, 103
- Sun Enterprise E10000, consola de administración, 29
- Sun Management Center, 37
- SunMC, *Ver* Sun Management Center
- SunPlex, *Ver* clúster
- SunPlex Manager, 37
- supervisión de las rutas de disco, 50
- supervisor de fallos, 72
- Supervisor de pertenencia al clúster, 40
- syncdir, opción de montaje, 49

T

- tiempo de CPU, 79
- tipos de recursos, 76
 - HASStoragePlus, 49
- topología de par en clúster, 35
- topología de pares en clúster, 30
- topología n+1 (estrella), 32
- topología n*n (escalable), 33
- topología par+n, 31
- topologías, 29, 34
 - n+1 (estrella), 32
 - n*n (escalable), 33
 - par en clúster, 35
 - par+n, 31
 - pares en clúster, 30
- transporte
 - cable, 101
 - ruta de acceso, 101

U

- UFS, 49
- unidad de CD-ROM, 26
- unidad de cinta, 26

V

- VERITAS Volume Manager, dispositivos multisistema, 25
- VxFS, 49