



# Sun Cluster の概念 (Solaris OS 版)



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

Part No: 820-6911-10  
2009 年 1 月、Revision A

Sun Microsystems, Inc. は、本書に記述されている技術に関する知的所有権を有しています。特に、この知的財産権はひとつかそれ以上の米国における特許、あるいは米国およびその他の国において申請中の特許を含んでいることがあります、それらに限定されるものではありません。

U.S. Government Rights – Commercial software. Government users are subject to the Sun Microsystems, Inc. standard license agreement and applicable provisions of the FAR and its supplements.

この配布には、第三者によって開発された素材を含んでいることがあります。

本製品の一部は、カリフォルニア大学からライセンスされている Berkeley BSD システムに基づいていることがあります。UNIX は、X/Open Company, Ltd. が独占的にライセンスしている米国ならびに他の国における登録商標です。

Sun、Sun Microsystems、Sun のロゴマーク、Solaris のロゴマーク、Java Coffee Cup のロゴマーク、docs.sun.com、OpenBoot、Solaris Volume Manager、StorEdge、Sun Fire、Java、および Solaris は、米国およびその他の国における米国 Sun Microsystems, Inc. (以下、米国 Sun Microsystems 社とします) またはその子会社の商標もしくは、登録商標です。すべての SPARC 商標は、米国 SPARC International, Inc. のライセンスを受けて使用している同社の米国およびその他の国における商標または登録商標です。SPARC 商標が付いた製品は、米国 Sun Microsystems 社が開発したアーキテクチャに基づくものです。

OPEN LOOK および Sun<sup>TM</sup> Graphical User Interface は、米国 Sun Microsystems 社が自社のユーザおよびライセンス実施権者向けに開発しました。米国 Sun Microsystems 社は、コンピュータ産業用のビジュアルまたはグラフィカルユーザインタフェースの概念の研究開発における米国 Xerox 社の先駆者としての成果を認めるものです。米国 Sun Microsystems 社は米国 Xerox 社から Xerox Graphical User Interface の非独占的ライセンスを取得しており、このライセンスは、OPEN LOOK のグラフィカル・ユーザインタフェースを実装するか、またはその他の方法で米国 Sun Microsystems 社との書面によるライセンス契約を遵守する、米国 Sun Microsystems 社のライセンス実施権者にも適用されます。

本書で言及されている製品や含まれている情報は、米国輸出規制法で規制されるものであり、その他の国の輸出入に関する法律の対象となることがあります。核、ミサイル、化学あるいは生物兵器、原子力の海洋輸送手段への使用は、直接および間接を問わず厳しく禁止されています。米国が禁輸の対象としている国や、限定はされませんが、取引禁止顧客や特別指定国民のリストを含む米国輸出排除リストで指定されているものへの輸出および再輸出は厳しく禁止されています。

本書は、「現状のまま」をベースとして提供され、商品性、特定目的への適合性または第三者の権利の非侵害の黙示の保証を含みそれに限定されない、明示的であるか黙示的であるかを問わない、なんらの保証も行われないものとします。

# 目次

---

|   |           |
|---|-----------|
| はじめに .....                                    | 7         |
| <b>1 基本知識と概要 .....</b>                        | <b>13</b> |
| Sun Cluster 環境の概要 .....                       | 13        |
| Sun Cluster ソフトウェアの3つの観点 .....                | 15        |
| ハードウェア保守担当者 .....                             | 16        |
| システム管理者 .....                                 | 17        |
| アプリケーション開発者 .....                             | 18        |
| Sun Cluster ソフトウェアの作業 .....                   | 19        |
| <b>2 重要な概念-ハードウェアサービスプロバイダ .....</b>          | <b>21</b> |
| Sun Cluster システムのハードウェアおよびソフトウェアコンポーネント ..... | 21        |
| クラスタノード .....                                 | 22        |
| クラスタハードウェアメンバー用のソフトウェアコンポーネント .....           | 23        |
| 多重ホストデバイス .....                               | 25        |
| 多重イニシエータ SCSI .....                           | 25        |
| ローカルディスク .....                                | 26        |
| リムーバブルメディア .....                              | 27        |
| クラスタインターコネクト .....                            | 27        |
| パブリックネットワークインタフェース .....                      | 28        |
| クライアントシステム .....                              | 28        |
| コンソールアクセスデバイス .....                           | 29        |
| 管理コンソール .....                                 | 30        |
| SPARC: Sun Cluster トポロジ .....                 | 31        |
| SPARC: クラスタペアトポロジ .....                       | 31        |
| SPARC: ペア +N トポロジ .....                       | 32        |
| SPARC: N+1 (星形) トポロジ .....                    | 33        |

|   |           |
|---|-----------|
| SPARC: N*N (スケラブル) トポロジ .....                     | 34        |
| SPARC: LDoms ゲストドメイン: ボックス内クラスタトポロジ .....         | 35        |
| SPARC: LDoms ゲストドメイン: 2つのホストにわたる単一クラスタトポロジ .....  | 37        |
| SPARC: LDoms ゲストドメイン: 2つのホストにわたる複数のクラスタトポロジ ..... | 38        |
| SPARC: LDoms ゲストドメイン: 冗長 I/O ドメイン .....           | 40        |
| x86: Sun Cluster トポロジ .....                       | 42        |
| x86: クラスタペアトポロジ .....                             | 42        |
| x86: N+1 (星形) トポロジ .....                          | 43        |
| <b>3 重要な概念 - システム管理者とアプリケーション開発者 .....</b>        | <b>45</b> |
| 管理インタフェース .....                                   | 46        |
| クラスタ内の時間 .....                                    | 46        |
| 高可用性フレームワーク .....                                 | 47        |
| ゾーンメンバーシップ .....                                  | 48        |
| クラスタメンバーシップモニター .....                             | 48        |
| フェイルファースト機構 .....                                 | 48        |
| クラスタ構成レポジトリ (CCR) .....                           | 50        |
| グローバルデバイス .....                                   | 50        |
| デバイス ID と DID 疑似ドライバ .....                        | 51        |
| デバイスグループ .....                                    | 51        |
| デバイスグループのフェイルオーバー .....                           | 52        |
| 多重ポートデバイスグループ .....                               | 53        |
| 広域名前空間 .....                                      | 55        |
| ローカル名前空間と広域名前空間の例 .....                           | 56        |
| クラスタファイルシステム .....                                | 56        |
| クラスタファイルシステムの用法 .....                             | 57        |
| HAStoragePlus リソースタイプ .....                       | 58        |
| syncdir マウントオプション .....                           | 59        |
| ディスクパスの監視 .....                                   | 60        |
| DPM の概要 .....                                     | 60        |
| ディスクパスの監視 .....                                   | 61        |
| 定足数と定足数デバイス .....                                 | 63        |
| 定足数投票数について .....                                  | 64        |
| 定足数の構成について .....                                  | 65        |

|   |     |
|---|-----|
| 定足数デバイス要件の順守 .....  | 65  |
| 定足数デバイスのベストプラクティスの順守 .....  | 66  |
| 推奨される定足数の構成 .....   | 67  |
| 変則的な定足数の構成 .....  | 69  |
| 望ましくない定足数の構成 .....  | 70  |
| データサービス .....   | 71  |
| データサービスメソッド .....   | 74  |
| フェイルオーバーデータサービス .....   | 75  |
| スケーラブルデータサービス .....   | 75  |
| 負荷均衡ポリシー .....  | 77  |
| フェイルバック設定 .....   | 79  |
| データサービス障害モニター .....   | 79  |
| 新しいデータサービスの開発 .....   | 80  |
| スケーラブルサービスの特徴 .....   | 80  |
| データサービス API と DSDL API .....  | 82  |
| クラスタインターコネクトによるデータサービストラフィックの送受信 .....  | 82  |
| リソース、リソースグループ、リソースタイプ .....   | 84  |
| リソースグループマネージャー (Resource Group Manager、RGM) .....                                   | 85  |
| リソースおよびリソースグループの状態と設定値 .....  | 85  |
| リソースとリソースグループプロパティ .....  | 87  |
| Solaris ゾーンをサポート .....  | 88  |
| RGM によるグローバルクラスタ非投票ノード (Solaris ゾーン) の直接サポート .....                                  | 88  |
| Sun Cluster ノード上の Solaris ゾーンを Sun Cluster HA for Solaris Containers を通してサポート ..... | 90  |
| サービス管理機能 .....  | 91  |
| システムリソースの使用状況 .....   | 92  |
| システムリソース監視 .....  | 93  |
| CPU の制御 .....   | 93  |
| システムリソース使用率の表示 .....  | 94  |
| データサービスプロジェクトの構成 .....  | 95  |
| プロジェクト構成に応じた要件の決定 .....   | 97  |
| プロセス当たりの仮想メモリー制限の設定 .....   | 98  |
| フェイルオーバーシナリオ .....  | 99  |
| パブリックネットワークアダプタと IP ネットワークマルチパス .....   | 105 |
| SPARC: 動的再構成のサポート .....   | 106 |
| SPARC: 動的再構成の概要 .....   | 107 |

---

|   |            |
|---|------------|
| SPARC: CPU デバイスに対する DR クラスタリング .....              | 107        |
| SPARC: メモリーに対する DR クラスタリング .....                  | 108        |
| SPARC: ディスクドライブとテープドライブに対する DR クラスタリング .....      | 108        |
| SPARC: 定足数デバイスに対する DR クラスタリング .....               | 108        |
| SPARC: クラスタインターコネクティブインタフェースに対する DR クラスタリング ..... | 109        |
| SPARC: パブリックネットワークインタフェースに対する DR クラスタリング .....    | 109        |
| <br>  |            |
| <b>4 よくある質問 .....</b>                             | <b>111</b> |
| 高可用性に関する FAQ .....                                | 111        |
| ファイルシステムに関する FAQ .....                            | 112        |
| ボリューム管理に関する FAQ .....                             | 114        |
| データサービスに関する FAQ .....                             | 114        |
| パブリックネットワークに関する FAQ .....                         | 115        |
| クラスタメンバーに関する FAQ .....                            | 116        |
| クラスタ記憶装置に関する FAQ .....                            | 117        |
| クラスタインターコネクティブに関する FAQ .....                      | 117        |
| クライアントシステムに関する FAQ .....                          | 118        |
| 管理コンソールに関する FAQ .....                             | 118        |
| 端末集配信装置とシステムサービスプロセッサに関する FAQ .....               | 119        |
| <br>  |            |
| 索引 .....  | 123        |

# はじめに

---

『Sun Cluster の概念 (Solaris OS 版)』には、SPARC® と x86 ベース両方のシステムの Sun™ Cluster 製品に関する概念的情報と参照情報が記載されています。

---

注 - この Sun Cluster リリースでは、SPARC および x86 系列のプロセッサアーキテクチャ (UltraSPARC、SPARC64、AMD64、および Intel 64) を使用するシステムをサポートします。このドキュメントでは、x86 とは 64 ビット x86 互換製品の広範囲なファミリーを指します。このドキュメントの情報では、特に明示されている場合以外はすべてのプラットフォームに関係します。

---

## 対象読者

このマニュアルは次の読者を対象としています。

- クラスタハードウェアを設置して保守を行う担当者
- Sun Cluster ソフトウェアをインストール、構成、管理するシステム管理者
- 現在 Sun Cluster 製品に含まれていないアプリケーション用のフェイルオーバーサービスやスケラブルサービスを開発するアプリケーション開発者

このマニュアルで説明されている概念を理解するには、Solaris オペレーティングシステムに精通し、Sun Cluster 製品とともに使用できるボリューム管理ソフトウェアに関する専門知識が必要です。

このマニュアルを読む前に、システムの必要条件を確認し、必要な装置とソフトウェアを購入しておく必要があります。『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』には、Sun Cluster ソフトウェアを計画、インストール、設定、および使用する方法が記載されています。

## 内容の紹介

『Sun Cluster の概念 (Solaris OS 版)』は、以下の章で構成されています。

第1章「基本知識と概要」では、Sun Cluster について知っておく必要がある全般的な概念の概要を説明します。

第2章「重要な概念 - ハードウェアサービスプロバイダ」では、ハードウェアサービスプロバイダが精通している必要がある概念を説明しています。これらの概念は、サービスプロバイダがハードウェアコンポーネント間の関係を理解するのに役立ちます。またこれらの概念は、サービスプロバイダとクラスタ管理者がクラスタソフトウェアおよびハードウェアをインストール、構成、管理する方法をよりよく理解することにも役立ちます。

第3章「重要な概念 - システム管理者とアプリケーション開発者」では、Sun Cluster API (Application Programming Interface) を使用するシステム管理者および開発者が知っておく必要がある概念を説明しています。開発者の方は、この API を使って、Web ブラウザやデータベースといった標準ユーザーアプリケーションを Sun Cluster 環境で動作する高可用性データサービスに変えることができます。

第4章「よくある質問」には、Sun Cluster 製品に関するよくある質問の回答が記載されています。

## 関連マニュアル

関連のある Sun Cluster のトピックについては、次の表に示したマニュアルを参照してください。すべての Sun Cluster マニュアルは、<http://docs.sun.com> で参照できます。

| 項目            | マニュアル  |
|---------------|--|
| 概要            | 『Sun Cluster の概要 (Solaris OS 版)』                                       |
|               | 『Sun Cluster 3.2 1/09 Documentation Center 』                           |
| 概念            | 『Sun Cluster の概念 (Solaris OS 版) 』                                      |
| ハードウェアの設計と管理  | 『Sun Cluster 3.1 - 3.2 Hardware Administration Manual for Solaris OS 』 |
|               | 各ハードウェア管理ガイド   |
| ソフトウェアのインストール | 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』                             |
|               | 『Sun Cluster クイックスタートガイド (Solaris OS 版)』                               |



| 項目                | マニュアル  |
|-------------------|--|
| データサービスのインストールと管理 | 『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』<br>各データサービスガイド  |
| データサービスの開発        | 『Sun Cluster データサービス開発ガイド (Solaris OS 版)』  |
| システム管理            | 『Sun Cluster のシステム管理 (Solaris OS 版)』<br>『Sun Cluster Quick Reference 』   |
| ソフトウェアアップグレード     | 『Sun Cluster Upgrade Guide for Solaris OS 』  |
| エラーメッセージ          | 『Sun Cluster Error Messages Guide for Solaris OS 』   |
| コマンドと関数のリファレンス    | 『Sun Cluster Reference Manual for Solaris OS 』<br>『Sun Cluster Data Services Reference Manual for Solaris OS 』<br>『Sun Cluster Quorum Server Reference Manual for Solaris OS 』 |

Sun Cluster ドキュメントの完全なリストについては、<http://wikis.sun.com/display/SunCluster/Home/> で Sun Cluster ソフトウェアの使用しているリリースのリリースノートを参照してください。

## 問い合わせについて

Sun Cluster ソフトウェアのインストールや使用に関して問題がある場合は、以下の情報をご用意の上、担当のサービスプロバイダにお問い合わせください。

- 名前と電子メールアドレス (利用している場合)
- 会社名、住所、および電話番号
- システムのモデルとシリアル番号
- オペレーティングシステムのバージョン番号 (例: Solaris 10 OS)
- Sun Cluster ソフトウェアのバージョン番号 (例: 3.2 1/09)

次のコマンドを使用し、システムに関して、サービスプロバイダに必要な情報を収集してください。

| コマンド                    | 機能                         |
|-------------------------|----------------------------|
| <code>prtconf -v</code> | システムメモリのサイズと周辺デバイス情報を表示します |
| <code>psrinfo -v</code> | プロセッサの情報を表示する              |
| <code>showrev -p</code> | インストールされているパッチを報告する        |

| コマンド                             | 機能                                  |
|----------------------------------|-------------------------------------|
| SPARC:prtdiag -v                 | システム診断情報を表示する                       |
| /usr/cluster/bin/clnode show-rev | Sun Cluster のリリースとパッケージバージョン情報を表示する |

上記の情報にあわせて、`/var/adm/messages` ファイルの内容もご購入先にお知らせください。

## マニュアル、サポート、およびトレーニング

Sun の Web サイトでは、次のサービスに関する情報も提供しています。

- マニュアル (<http://jp.sun.com/documentation/>)
- サポート (<http://jp.sun.com/support/>)
- トレーニング (<http://jp.sun.com/training/>)

## 表記上の規則

このマニュアルでは、次のような字体や記号を特別な意味を持つものとして使用します。

表 P-1 表記上の規則

| 字体または記号          | 意味  | 例   |
|------------------|---|---|
| AaBbCc123        | コマンド名、ファイル名、ディレクトリ名、画面上のコンピュータ出力、コード例を示します。 | .login ファイルを編集します。<br>ls -a を使用してすべてのファイルを表示します。<br><br>system% |
| <b>AaBbCc123</b> | ユーザーが入力する文字を、画面上のコンピュータ出力と区別して示します。         | system% <b>su</b><br>password:                                  |
| <i>AaBbCc123</i> | 変数を示します。実際に使用する特定の名前または値で置き換えます。            | ファイルを削除するには、 <code>rm filename</code> と入力します。                   |
| 『 』              | 参照する書名を示します。                                | 『コードマネージャ・ユーザーズガイド』を参照してください。                                   |

表 P-1 表記上の規則 (続き)

| 字体または記号 | 意味                                     | 例  |
|---------|--|--|
| 「」      | 参照する章、節、ボタンやメニュー名、強調する単語を示します。         | 第5章「衝突の回避」を参照してください。<br><br>この操作ができるのは、「スーパーユーザー」だけです。 |
| \       | 枠で囲まれたコード例で、テキストがページ行幅を超える場合に、継続を示します。 | sun% <b>grep</b> '^#define \<br><br>XV_VERSION_STRING' |

コード例は次のように表示されます。

- C シェル

```
machine_name% command y|n [filename]
```

- C シェルのスーパーユーザー

```
machine_name# command y|n [filename]
```

- Bourne シェルおよび Korn シェル

```
$ command y|n [filename]
```

- Bourne シェルおよび Korn シェルのスーパーユーザー

```
# command y|n [filename]
```

[ ] は省略可能な項目を示します。上記の例は、*filename* は省略してもよいことを示しています。

| は区切り文字 (セパレータ) です。この文字で分割されている引数のうち 1 つだけを指定します。

キーボードのキー名は英文で、頭文字を大文字で示します (例: Shift キーを押します)。ただし、キーボードによっては Enter キーが Return キーの動作をします。

ダッシュ (-) は 2 つのキーを同時に押すことを示します。たとえば、Ctrl-D は Control キーを押したまま D キーを押すことを意味します。



# 基本知識と概要

---

Sun Cluster 製品はハードウェアとソフトウェアが統合されたソリューションであり、高度な可用性とスケーラビリティを備えたサービスを作成するために使用されます。『Sun Cluster の概念 (Solaris OS 版)』では、Sun Cluster 製品をより深く理解するために必要な概念の情報を説明します。このマニュアルは、Sun Cluster のほかのマニュアルと合わせて、Sun Cluster ソフトウェアの全体を説明するものです。

この章では、Sun Cluster 製品の根底にある一般的な概念の概要を説明します。

この章では次の内容を示します。

- Sun Cluster ソフトウェアの基本知識と概要
- 各ユーザーから見た Sun Cluster
- Sun Cluster ソフトウェアを使用する前に理解しておく必要がある重要な概念の明示
- 重要な概念に関連する手順と情報を記載した Sun Cluster のマニュアル
- クラスタに関連する作業と、これらの作業手順が記載されたマニュアル

この章で説明する内容は次のとおりです。

- 13 ページの「[Sun Cluster 環境の概要](#)」
- 15 ページの「[Sun Cluster ソフトウェアの 3 つの観点](#)」
- 19 ページの「[Sun Cluster ソフトウェアの作業](#)」

## Sun Cluster 環境の概要

Sun Cluster 環境は、Solaris オペレーティングシステムをクラスタオペレーティングシステムに拡張するものです。「クラスタ」は、1 つまたは複数のノードの集合です。各ノードは、この集合に排他的に属しています。Solaris 10 OS 上で動作するクラスタには、「グローバルクラスタ」と「ゾーンクラスタ」の 2 種類があります。

Solaris 10 OS より前にリリースされた Solaris OS 上で動作するクラスタの場合、ノードはクラスタメンバーシップに参加する「物理的なマシン」であり、定数数デバイ

スではありません。Solaris 10 OS 上で動作するクラスタでは、ノードの概念が変更されました。この環境でのノードは、クラスタに関連付けられている Solaris ゾーンです。また、「Solaris ホスト」(単純に「ホスト」とも呼ばれます)は、Solaris OS および個別のプロセスが実行される、次のハードウェア構成またはソフトウェア構成のいずれかです。

- 仮想マシンで構成されていない、またはハードウェアドメインとして構成されていない「ベアメタル」物理マシン
- Sun Logical Domains (LDoms) のゲストドメイン
- Sun Logical Domains (LDoms) の I/O ドメイン
- ハードウェアドメイン

これらのプロセスは、相互にやりとりすることによって、ユーザーに提供するアプリケーション、システムリソース、データを(ネットワーククライアントにとって)1つのシステムのように形成します。

Solaris 10 環境のグローバルクラスタは、1つまたは複数のグローバルクラスタ投票ノード、および任意で0または1つ以上のグローバルクラスタ非投票ノードだけで構成されるクラスタの一種です。

---

注- グローバルクラスタには、オプションで solaris8、solaris9、lx(linux)、またはネイティブブランドの非大域ゾーンを含めることができます。これらはノードではなく、高可用性のコンテナ(リソース)です。

---

グローバルクラスタ投票ノードとは、グローバルクラスタ内のネイティブブランドの大域ゾーンで、定足数投票、つまりクラスタのメンバーシップ投票の、総数の票を構成します。この総数により、そのクラスタが処理を継続するのに十分な票を持っているかどうかが決まります。グローバルクラスタ非投票ノードとは、グローバルクラスタ内のネイティブブランドの非大域ゾーンで、定足数投票、つまりクラスタのメンバーシップ投票の、総数の票を構成しません。

Solaris 10 環境では、ゾーンクラスタは、1つまたは複数のクラスタブランドの投票ノードのみから構成されるクラスタです。ゾーンクラスタは、グローバルクラスタに依存しており、したがって、グローバルクラスタを必要とします。グローバルクラスタはゾーンクラスタを含みません。ゾーンクラスタを構成するには、グローバルクラスタが必要です。ゾーンクラスタは1つのマシン上に最大で1つのゾーンクラスタノードを持ちます。

---

注-ゾーンクラスタノードは、同一マシン上のグローバルクラスタ投票ノードが処理を継続している間にかぎり、処理を継続できます。マシン上でグローバルクラスタ投票ノードに障害が発生すると、同じマシンのすべてのゾーンクラスタノードにも障害が発生します。

---

クラスタには、従来の単一サーバーシステムと比較した場合、いくつかの利点があります。これらの利点には、フェイルオーバーサービスとスケラブルサービスのサポート、モジュールの成長に対応できる容量、従来のハードウェアフォルトトレラントシステムよりも低価格の製品といったものがあります。

次に、Sun Cluster ソフトウェアの目的を示します。

- ソフトウェアまたはハードウェアの障害が原因のシステム停止時間を短縮するか完全になくします。
- 単一サーバーシステムを停止させるような障害が発生しても、エンドユーザーへのデータとアプリケーションの可用性を保証します。
- クラスタにノードを追加し、追加したプロセッサに応じたサービスを提供できるようにすることで、アプリケーションのスループットを向上させます。
- クラスタ全体を停止しなくても保守を実行できるようにすることで、システムの可用性を強化します。

フォルトトレラント機能と高可用性についての詳細は、『[Sun Cluster の概要 \(Solaris OS 版\)](#)』の「[Sun Cluster によるアプリケーションの可用性の向上](#)」を参照してください。

高可用性の FAQ については、[111 ページ](#)の「[高可用性に関する FAQ](#)」を参照してください。

## Sun Cluster ソフトウェアの3つの観点

この節では、Sun Cluster ソフトウェアのユーザーを3種類に分け、各ユーザーに関連する概念とマニュアルについて説明します。

各ユーザーは次のとおりです。

- ハードウェア保守担当者
- システム管理者
- アプリケーション開発者

## ハードウェア保守担当者

ハードウェア保守担当者にとって、Sun Cluster ソフトウェアは、サーバー、ネットワーク、および記憶装置を含む市販のハードウェアの集合に見えます。これらのコンポーネントは、すべてのコンポーネントにバックアップがあり、単一の障害によってシステム全体が停止しないように配線されています。

### 重要な概念 (ハードウェア保守担当者)

ハードウェア保守担当者は、クラスタに関する次の概念を理解する必要があります。

- クラスタハードウェアの構成と配線
- 設置と保守 (追加、取り外し、交換)
  - ネットワークインタフェースコンポーネント (アダプタ、接続点、ケーブル)
  - ディスクインタフェースカード
  - ディスクアレイ
  - ディスクドライブ
  - 管理コンソールとコンソールアクセスデバイス
- 管理コンソールとコンソールアクセスデバイスの設定

### 参照箇所 (ハードウェア保守担当者)

次の項には、前述の重要な概念に関連する説明が記載されています。

- 22 ページの「クラスタノード」
- 25 ページの「多重ホストデバイス」
- 26 ページの「ローカルディスク」
- 27 ページの「クラスタインターコネクト」
- 28 ページの「パブリックネットワークインタフェース」
- 28 ページの「クライアントシステム」
- 30 ページの「管理コンソール」
- 29 ページの「コンソールアクセスデバイス」
- 31 ページの「SPARC: クラスタペアトポロジ」
- 33 ページの「SPARC: N+1 (星形) トポロジ」

### Sun Cluster の関連マニュアル (ハードウェア保守担当者)

『Sun Cluster 3.1 - 3.2 Hardware Administration Manual for Solaris OS』には、ハードウェアサービスの概念に関連する手順と情報が記載されています。



## システム管理者

システム管理者にとって、Sun Cluster 製品は、記憶装置を共有する Solaris ホストの集合です。

システム管理者は、次の作業を行うソフトウェアを扱います。

- クラスタ内の Solaris ホスト間のコネクティビティーを監視するための、Solaris ソフトウェアに統合された専用のクラスタソフトウェア
- クラスタノードで実行されるユーザーアプリケーションプログラムの状態を監視するための専用のソフトウェア
- ディスクを設定して管理するためのボリュームマネージャー
- 直接ディスクに接続されていない Solaris ホストも含め、すべての Solaris ホストが、すべての記憶装置にアクセスできるようにするための専用のクラスタソフトウェア
- ファイルがすべての Solaris ホストに対してローカルに接続されているように表示するための専用のクラスタソフトウェア

### 重要な概念(システム管理者)

システム管理者は、次の概念とプロセスについて理解する必要があります。

- ハードウェアとソフトウェアの間の対話
- クラスタをインストールして構成する方法の一般的な流れ
  - Solaris オペレーティングシステムのインストール
  - Sun Cluster ソフトウェアのインストールと構成
  - ボリュームマネージャーのインストールと構成
  - クラスタを動作可能状態にするためのアプリケーションソフトウェアのインストールと構成
  - Sun Cluster データサービスソフトウェアのインストールと構成
- クラスタハードウェアとソフトウェアのコンポーネントを追加、削除、交換、およびサービス提供するためのクラスタ管理手順
- パフォーマンスを向上させるための構成の変更方法

### 参照個所(システム管理者)

次の項には、前述の重要な概念に関連する説明が記載されています。

- 46 ページの「管理インタフェース」
- 46 ページの「クラスタ内の時間」
- 47 ページの「高可用性フレームワーク」
- 50 ページの「グローバルデバイス」

- 51 ページの「デバイスグループ」
- 55 ページの「広域名前空間」
- 56 ページの「クラスタファイルシステム」
- 60 ページの「ディスクパスの監視」
- 71 ページの「データサービス」

## システム管理者向けの Sun Cluster のマニュアル

次の Sun Cluster のマニュアルには、システム管理者の概念に関連する手順と情報が記載されています。

- 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』
- 『Sun Cluster のシステム管理 (Solaris OS 版)』
- 『Sun Cluster Error Messages Guide for Solaris OS』
- 『Sun Cluster リリースノートご使用にあたって (Solaris OS 版)』

## アプリケーション開発者

Sun Cluster ソフトウェアは、Oracle、NFS、DNS、Sun Java System Web Server、Apache Web Server (SPARC ベースシステム上)、Sun Java System Directory Server などのアプリケーションに対応する「データサービス」を提供します。データサービスを作成するには、既成のアプリケーションを Sun Cluster ソフトウェアの制御下で動作するように設定する必要があります。Sun Cluster ソフトウェアは、このようなアプリケーションの起動、停止、および監視を行う構成ファイルと管理メソッドを提供します。新しいフェイルオーバーサービスまたはスケラブルサービスを作成する必要がある場合は、Sun Cluster Application Programming Interface (API) と Data Service Enabling Technologies API (DSET API) を使用して、そのアプリケーションがクラスタ上でデータサービスとして実行するために必要な構成ファイルと管理メソッドを開発します。

### 重要な概念 (アプリケーション開発者)

アプリケーション開発者は、次の概念について理解する必要があります。

- 開発するアプリケーションの特性。その特性に基づいて、アプリケーションをフェイルオーバーまたはスケラブルデータサービスとして実行できるかどうかを判断する必要があります。
- Sun Cluster API、DSET API、汎用データサービス。開発者は、各自のアプリケーションをクラスタ環境に合わせて構成するプログラムまたはスクリプトを記述するために、どのツールがもっとも適しているかを判断する必要があります。

## 参照箇所(アプリケーション開発者)

次の項には、前述の重要な概念に関連する説明が記載されています。

- 71 ページの「データサービス」
- 84 ページの「リソース、リソースグループ、リソースタイプ」
- 第4章「よくある質問」

## アプリケーション開発者向けの Sun Cluster のマニュアル

次の Sun Cluster のマニュアルには、アプリケーション開発者の概念に関連する手順と情報が記載されています。

- 『Sun Cluster データサービス開発ガイド (Solaris OS 版)』
- 『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』

# Sun Cluster ソフトウェアの作業

すべての Sun Cluster ソフトウェアの作業には、ある程度の概念的な背景知識が必要です。次の表は、作業と作業手順が記載されたマニュアルを示したものです。このマニュアルの概念に関する章では、各概念がこれらの作業とどのように対応するかを説明します。

表 1-1 作業マップ: ユーザーの作業と参照するマニュアル

| 作業  | 参照先   |
|---|---|
| クラスタハードウェアの設置                               | 『Sun Cluster 3.1 - 3.2 Hardware Administration Manual for Solaris OS』 |
| クラスタへの Solaris ソフトウェアのインストール                | 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』                            |
| SPARC: Sun™ Management Center ソフトウェアのインストール | 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』                            |
| Sun Cluster ソフトウェアのインストールと構成                | 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』                            |
| ボリュームマネージャーのインストールと構成                       | 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』<br>各ボリュームマネージャーのマニュアル      |
| Sun Cluster データサービスのインストールと構成               | 『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』                            |
| クラスタハードウェアの保守                               | 『Sun Cluster 3.1 - 3.2 Hardware Administration Manual for Solaris OS』 |

表 1-1 作業マップ: ユーザーの作業と参照するマニュアル (続き)

| 作業                    | 参照先  |
|-----------------------|--|
| Sun Cluster ソフトウェアの管理 | 『Sun Cluster のシステム管理 (Solaris OS 版)』                     |
| ボリュームマネージャーの管理        | 『Sun Cluster のシステム管理 (Solaris OS 版)』 およびボリューム管理に関するマニュアル |
| アプリケーションソフトウェアの管理     | 各アプリケーションのマニュアル  |
| 問題の識別と対処方法            | 『Sun Cluster Error Messages Guide for Solaris OS 』       |
| 新しいデータサービスの作成         | 『Sun Cluster データサービス開発ガイド (Solaris OS 版)』                |

## 重要な概念 - ハードウェアサービスプロバイダ

---

この章では、Sun Cluster 構成のハードウェアコンポーネントに関連する重要な概念について説明します。

この章の内容は次のとおりです。

- 21 ページの「Sun Cluster システムのハードウェアおよびソフトウェアコンポーネント」
- 31 ページの「SPARC: Sun Cluster トポロジ」
- 42 ページの「x86: Sun Cluster トポロジ」

## Sun Cluster システムのハードウェアおよびソフトウェアコンポーネント

ここで示す情報は、主にハードウェアサービスプロバイダを対象としています。これらの概念は、サービスプロバイダが、クラスタハードウェアの設置、構成、またはサービスを提供する前に、ハードウェアコンポーネント間の関係を理解するのに役立ちます。またこれらの情報は、クラスタシステムの管理者にとっても、クラスタソフトウェアをインストール、構成、管理するための予備知識として役立ちます。

クラスタは、次のようなハードウェアコンポーネントで構成されます。

- ローカルディスク (非共有) を備えた Solaris ホスト
- 多重ホスト記憶装置 (Solaris ホスト間で共有されるディスク)
- リムーバブルメディア (テープ、CD-ROM)
- クラスタインターコネクタ
- パブリックネットワークインタフェース
- クライアントシステム
- 管理コンソール
- コンソールアクセスデバイス

Sun Cluster ソフトウェアを使用すると、これらのコンポーネントを各種の構成に組み合わせることができます。これらの構成については、次の節で説明します。

- 31 ページの「SPARC: Sun Cluster トポロジ」
- 42 ページの「x86: Sun Cluster トポロジ」

2 ホストクラスタの構成例については、『Sun Cluster の概要 (Solaris OS 版)』の「Sun Cluster のハードウェア環境」を参照してください。

## クラスタノード

Solaris 10 OS より前にリリースされた Solaris OS 上で動作するクラスタの場合、ノードはクラスタメンバーシップに参加する「物理的なマシン」であり、定数デバイスではありません。Solaris 10 OS 上で動作するクラスタでは、ノードの概念が変更されました。この環境でのノードは、クラスタに関連付けられている Solaris ゾーンです。また、「Solaris ホスト」(単純に「ホスト」とも呼ばれます)は、Solaris OS および個別のプロセスが実行される、次のハードウェア構成またはソフトウェア構成のいずれかです。

- 仮想マシンで構成されていない、またはハードウェアドメインとして構成されていない「ベアメタル」物理マシン
- Sun Logical Domains (LDoms) のゲストドメイン
- Sun Logical Domains (LDoms) の I/O ドメイン
- ハードウェアドメイン

プラットフォームに応じて、Sun Cluster ソフトウェアは次の構成をサポートします。

- SPARC: Sun Cluster ソフトウェアは、1つのクラスタで1つから16までの Solaris ホストをサポートします。ハードウェア構成によっては、SPARC ベースのシステムから成るクラスタで構成できるホストの最大数に制限が追加されます。サポートされる構成については、31 ページの「SPARC: Sun Cluster トポロジ」を参照してください。
- x86: Sun Cluster ソフトウェアは、1つのクラスタで1つから8つまでの Solaris ホストをサポートします。ハードウェア構成によっては、x86 ベースのシステムから成るクラスタで構成できるホストの最大数に制限が追加されます。サポートされる構成については、42 ページの「x86: Sun Cluster トポロジ」を参照してください。

一般的に Solaris ホストは、1つまたは複数の多重ホストデバイスに接続されます。多重ホストデバイスに接続されていないホストは、クラスタファイルシステムを使用して多重ホストデバイスにアクセスします。たとえば、スケラブルサービスを1つ構成することで、ホストが多重ホストデバイスに直接接続されていなくてもサービスを提供することができます。

さらに、パラレルデータベース構成では、複数のホストがすべてのディスクへの同時アクセスを共有します。

- ディスクへの同時アクセスについては、25 ページの「多重ホストデバイス」を参照してください。
- パラレルデータベース構成についての詳細は、31 ページの「SPARC: クラスタペアトポロジ」と42 ページの「x86: クラスタペアトポロジ」を参照してください。

クラスタ内のノードはすべて、共通の名前(クラスタ名)によってグループ化されます。この名前は、クラスタのアクセスと管理に使用されます。

パブリックネットワークアダプタは、ホストとパブリックネットワークを接続して、クラスタへのクライアントアクセスを可能にします。

クラスタメンバーは、1つまたは複数の物理的に独立したネットワークを介して、クラスタ内のほかのホストと通信します。物理的に独立したネットワークの集合は、クラスタインターコネクトと呼ばれます。

クラスタ内のすべてのノードは、別のノードがいつクラスタに結合されたか、またはクラスタから切り離されたかを認識します。さらに、クラスタ内のすべてのノードは、ほかのクラスタノードで実行されているリソースだけでなく、ローカルに実行されているリソースも認識します。

同じクラスタ内の各ホストの処理、メモリー、および入出力機能が同等で、パフォーマンスを著しく低下させることなく処理を継続できることを確認してください。フェイルオーバーの可能性があるため、すべてのホストには、バックアップまたは二次ホストとしてすべてのホストの作業負荷を引き受けるのに十分な予備容量が必要です。

各ホストは、独自のルート (/) ファイルシステムを起動します。

## クラスタハードウェアメンバー用のソフトウェアコンポーネント

Solaris ホストがクラスタメンバーとして動作するためには、ホストに次のソフトウェアがインストールされていなければなりません。

- Solaris オペレーティングシステム
- Sun Cluster ソフトウェア
- データサービスアプリケーション
- ボリューム管理 (Solaris Volume Manager™ または Veritas Volume Manager)

例外として、複数のディスクの冗長配列 (RAID) を使用する構成があります。この構成には、通常、Solaris Volume Manager や Veritas Volume Manager などのボリュームマネージャーは必要ありません。

- Solaris オペレーティングシステム、Sun Cluster、およびボリュームマネージャーのインストール方法については、『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』を参照してください。
- データサービスのインストールおよび構成については、『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』を参照してください。
- 前述のソフトウェアコンポーネントの概念については、第3章「重要な概念 - システム管理者とアプリケーション開発者」を参照してください。

次の図は、Sun Cluster 環境を構成するソフトウェアコンポーネントとその関係の概要を示しています。

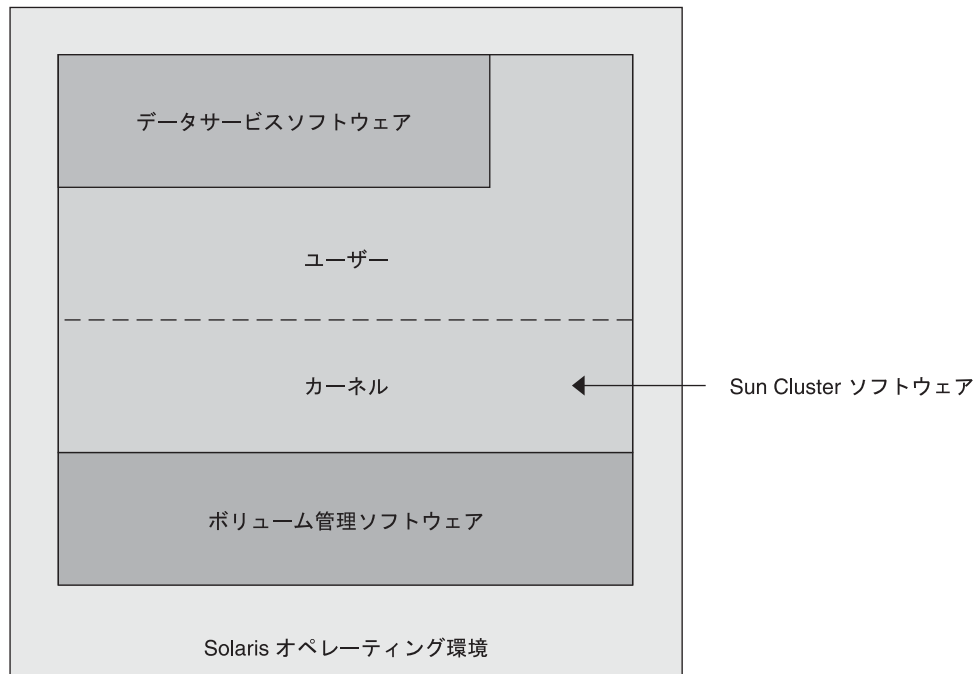


図 2-1 Sun Cluster ソフトウェアコンポーネントとその関係の概要

クラスタメンバーの FAQ については、第4章「よくある質問」を参照してください。



## 多重ホストデバイス

多重ホストデバイスとは、一度に複数の Solaris ホストに接続できるディスクのことです。Sun Cluster 環境では、多重ホスト記憶装置によってディスクの可用性を強化できます。2 ホストクラスタでは、Sun Cluster ソフトウェアは定足数を確立するために多重ホスト記憶装置を必要とします。3 ホストより大きなクラスタでは、定足数デバイスを必要としません。定足数についての詳細は、63 ページの「定足数と定足数デバイス」を参照してください。

多重ホストデバイスには、次の特徴があります。

- 単一ホスト障害への耐性。
- アプリケーションデータ、アプリケーションバイナリ、および構成ファイルを格納する機能。
- ホスト障害からの保護。クライアントがあるホストを介するデータを要求して、そのホストに障害が発生した場合、これらの要求は、同じディスクに直接接続されている別のホストを使用するようにスイッチオーバーされます。
- ディスクを「マスター」する主ホストを介する広域的なアクセス、またはローカルパスを介する直接同時アクセス。現在、直接同時アクセスを使用するアプリケーションは Oracle Real Application Clusters Guard だけです。

ボリュームマネージャーは、ミラー化された構成または RAID-5 構成を提供することによって、多重ホストデバイスのデータ冗長性を実現します。現在、Sun Cluster はボリュームマネージャーとして Solaris Volume Manager および Veritas Volume Manager をサポートし、また、いくつかのハードウェア RAID プラットフォームでは RDAC RAID-5 ハードウェアコントローラをサポートします。

多重ホストデバイスをミラー化したディスクやストライプ化したディスクと組み合わせると、ホストの障害や個々のディスクの障害から保護できます。

多重ホスト記憶装置の FAQ については、第 4 章「よくある質問」を参照してください。

## 多重イニシエータ SCSI

この項は、多重ホストデバイスに使用されるファイバチャネル記憶装置ではなく、SCSI 記憶装置にのみ適用されます。

クラスタ化されていないスタンドアロンホストでは、ホストが、このホストを特定の SCSI バスに接続する SCSI ホストアダプタ回路によって、SCSI バスのアクティビティを制御します。この SCSI ホストアダプタ回路は、SCSI イニシエータと呼ばれます。この回路は、この SCSI バスに対するすべてのバスアクティビティを開始します。Sun システムの SCSI ホストアダプタのデフォルト SCSI アドレスは 7 です。

クラスタ構成では、多重ホストデバイスを使用し、複数のホスト間で記憶装置を共有します。クラスタ記憶装置が SCSI デバイスまたは Differential SCSI デバイスで構成される場合、その構成のことを「多重イニシエータ SCSI」と呼びます。この用語が示すように、複数の SCSI イニシエータが SCSI バスに存在します。

SCSI 仕様では、SCSI バス上のデバイスごとに一意の SCSI アドレスが必要 (ホストアダプタも SCSI バス上のデバイス) です。多重イニシエータ環境では、デフォルトのハードウェア構成は、すべての SCSI ホストアダプタがデフォルトの 7 になっているので、衝突が生じます。

この衝突を解決するには、各 SCSI バスで、SCSI アドレスが 7 の SCSI ホストアダプタを 1 つ残し、ほかのホストアダプタには、未使用の SCSI アドレスを設定します。これらの未使用の SCSI アドレスには、現在未使用のアドレスと最終的に未使用となるアドレスの両方を含めるべきです。将来未使用となるアドレスの例としては、新しいドライブを空のドライブスロットに設置することによる記憶装置の追加がありません。

ほとんどの構成では、二次ホストアダプタに使用できる SCSI アドレスは 6 です。

これらのホストアダプタ用に選択された SCSI アドレスを変更するには、次のツールのいずれかを使用して、`scsi-initiator-id` プロパティを設定します。

- `eeeprom(1M)`
- SPARC ベースシステム上の OpenBoot™ PROM
- x86 ベースのシステムで BIOS のブート後に任意で実行する SCSI ユーティリティ

このプロパティは 1 つのホストに対して、広域的にまたはホストアダプタごとに設定できます。SCSI ホストアダプタごとに一意の `scsi-initiator-id` を設定する手順は、『Sun Cluster 3.1 - 3.2 With SCSI JBOD Storage Device Manual for Solaris OS』に記載されています。

## ローカルディスク

ローカルディスクとは、単一の Solaris ホストにのみ接続されたディスクを表します。したがって、ローカルディスクはホストの障害から保護されません。つまり、可用性が低いということです。ただし、ローカルディスクを含むすべてのディスクは広域の名前空間に含まれ、広域デバイスとして構成されています。したがって、ディスク自体はすべてのクラスタホストから参照できます。

ローカルディスク上のファイルシステムをほかのホストから使用できるようにするには、それらのファイルシステムを広域マウントポイントに置きます。これらの広域ファイルシステムのいずれかがマウントされているホストに障害が生じると、すべてのホストがそのファイルシステムにアクセスできなくなります。ボリュームマネージャーを使用すると、これらのディスクがミラー化されるため、これらのファイルシステムに障害が発生してもアクセス不能になることはありません。ただし、ホスト障害をボリュームマネージャーで保護することはできません。

広域デバイスについての詳細は、50 ページの「グローバルデバイス」を参照してください。

## リムーバブルメディア

クラスタでは、テープドライブや CD-ROM ドライブなどのリムーバブルメディアがサポートされています。通常、これらのデバイスは、クラスタ化していない環境と同じ方法でインストール、構成し、サービスを提供できます。これらのデバイスは、Sun Cluster で広域デバイスとして構成されるため、クラスタ内の任意のノードから各デバイスにアクセスできます。リムーバブルメディアのインストールと構成については、『Sun Cluster 3.1 - 3.2 Hardware Administration Manual for Solaris OS』を参照してください。

広域デバイスについての詳細は、50 ページの「グローバルデバイス」を参照してください。

## クラスタインターコネクト

「クラスタインターコネクト」は、クラスタ内の Solaris ホスト間のクラスタプライベート通信とデータサービス通信の転送に使用される物理的な装置構成です。インターコネクトは、クラスタプライベート通信で拡張使用されるため、パフォーマンスが制限される可能性があります。

クラスタ内のホストだけがクラスタインターコネクトに接続できます。Sun Cluster セキュリティモデルは、クラスタホストだけがクラスタインターコネクトに物理的にアクセスできるものと想定しています。

1つのクラスタでは、1つから6つまでのクラスタインターコネクトを設定できます。クラスタインターコネクトを1つだけ使用すると、プライベートインターコネクトに使用されるアダプタポートの数が減り、同時に冗長性がなくなり、可用性が低くなります。また、1つのインターコネクトに障害が発生すると、クラスタについて自動回復を実行しなければならないリスクが高くなります。可能な限り、クラスタインターコネクトは2つ以上インストールしてください。冗長性とスケーラビリティが提供されるので、シングルポイント障害が回避されて可用性も高くなります。

クラスタインターコネクトは、アダプタ、接続点、およびケーブルの3つのハードウェアコンポーネントで構成されます。次に、これらの各ハードウェアコンポーネントについて説明します。

- アダプタ - 個々のクラスタホストに存在するネットワークインタフェースカード。アダプタの名前は、デバイス名と物理ユニット番号で構成されます (qfe2 など)。一部のアダプタには物理ネットワーク接続が1つしかありませんが、qfe カードのように複数の物理接続を持つものもあります。また、ネットワークインタフェースと記憶装置インタフェースの両方を持つものもあります。

複数のインタフェースを持つネットワークアダプタは、アダプタ全体に障害が生じると、単一地点による障害の原因となる可能性があります。可用性を最適にするには、2つのホスト間の唯一のパスが単一のネットワークアダプタに依存しないように、クラスタを設定してください。

- 接続点 - クラスタホストの外部に存在するスイッチ。接続点は、パススルーおよび切り換え機能を実行して、3つ以上のホストに接続できるようにします。2ホストクラスタでは、各ホストの冗長アダプタに接続された冗長物理ケーブルによって、ホストを相互に直接接続できるため、接続点は必要ありません。3ホスト以上の構成では、通常は接続点が必要です。
- ケーブル - 2つのネットワークアダプタ間、アダプタと接続点の間に設置する物理接続。

クラスタインターコネクットの FAQ については、[第4章「よくある質問」](#)を参照してください。

## パブリックネットワークインタフェース

クライアントは、パブリックネットワークインタフェースを介してクラスタに接続します。各ネットワークアダプタカードは、カードに複数のハードウェアインタフェースがあるかどうかによって、1つまたは複数のパブリックネットワークに接続できます。

複数のパブリックネットワークインタフェースカードを持つ Solaris ホストをクラスタに設定することによって、次の機能を実行できます。

- 複数のカードをアクティブにするよう構成する。
- 相互のフェイルオーバーバックアップとする。

いずれかのアダプタに障害が発生すると、IP ネットワークマルチパスソフトウェアが呼び出され、障害のあるインタフェースが同じグループの別のアダプタにフェイルオーバーされます。

パブリックネットワークインタフェースのクラスタ化に関連する特殊なハードウェアについての特記事項はありません。

パブリックネットワークの FAQ については、[第4章「よくある質問」](#)を参照してください。

## クライアントシステム

クライアントシステムには、パブリックネットワークによってクラスタにアクセスするマシンやほかのホストが含まれます。クライアント側プログラムは、クラスタ上で動作しているサーバー側アプリケーションが提供するデータやサービスを使用します。

クライアントシステムの可用性は高くありません。クラスタ上のデータとアプリケーションは、高い可用性を備えています。

クライアントシステムの FAQ については、第 4 章「よくある質問」を参照してください。

## コンソールアクセスデバイス

クラスタ内のすべての Solaris ホストにはコンソールアクセスが必要です。

コンソールアクセスを取得するには、次のうちの 1 つのデバイスを使用します。

- クラスタハードウェアとともに購入した端末集配信装置
- Sun Enterprise E10000 サーバーのシステムサービスプロセッサ (System Service Processor、SSP) (SPARC ベースクラスタの場合)
- Sun Fire™ サーバーのシステムコントローラ (同じく SPARC ベースクラスタの場合)
- 各ホストの ttya にアクセスできる別のデバイス

サポートされている唯一の端末集配信装置は、Sun から提供されています。サポートされている Sun の端末集配信装置の使用は任意です。端末集配信装置を使用すると、TCP/IP ネットワークを使用して、各ホストの `/dev/console` にアクセスできます。この結果、ネットワークの任意の場所にあるリモートマシンから、各ホストにコンソールレベルでアクセスできます。

システムサービスプロセッサ (System Service Processor、SSP) は、Sun Enterprise E10000 サーバーへのコンソールアクセスを提供します。SSP とは、Sun Enterprise E10000 サーバーをサポートするように構成された Ethernet ネットワーク上のマシンのプロセッサカードのことです。SSP は、Sun Enterprise E10000 サーバーの管理コンソールです。Sun Enterprise E10000 サーバーのネットワークコンソール機能を使用すると、ネットワーク上のすべてのマシンからホストコンソールセッションを開くことができます。

これ以外のコンソールアクセス方式には、ほかの端末集配信装置、別ホストおよびダム端末からの tip シリアルポートアクセスがあります。



注意-基本サーバプラットフォームでキーボードまたはモニターがサポートされている場合、クラスタホストにキーボードまたはモニターを接続できます。ただし、このキーボードまたはモニターはコンソールデバイスとして使用できません。コンソールはシリアルポートにリダイレクトする必要があります。マシンによっては、適切な OpenBoot PROM パラメータを設定して、システムサービスプロセッサ (System Service Processor、SSP) およびリモートシステム制御 (Remote System Control、RSC) にコンソールをリダイレクトする必要があります。

## 管理コンソール

管理コンソールと呼ばれる専用のマシンを使用して動作中のクラスタを管理できます。通常は、Cluster Control Panel (CCP) や Sun Management Center 製品の Sun Cluster モジュール (SPARC ベースクラスタのみ) などの管理ツールソフトウェアを管理コンソールにインストールして実行します。CCP で `cconsole` を使用すると、一度に複数のホストコンソールに接続できます。CCP の使用法についての詳細は、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の第 1 章「[Sun Cluster の管理の概要](#)」を参照してください。

管理コンソールはクラスタホストではありません。管理コンソールは、パブリックネットワークを介して、または任意でネットワークベースの端末集配信装置を経由して、クラスタホストへのリモートアクセスに使用します。

クラスタが Sun Enterprise E10000 プラットフォームで構成されている場合は、次の作業を行います。

- 管理コンソールから SSP にログインする。
- `netcon` コマンドを使用して接続する。

通常、ホストはモニターなしで構成します。そして、管理コンソールから `telnet` セッションを使用して、ホストのコンソールにアクセスします。管理コンソールは端末集配信装置に接続され、端末集配信装置から当該ホストのシリアルポートに接続されます。Sun Enterprise E1000 サーバーの場合は、システムサービスプロセッサから接続します。詳細は、[29 ページ](#)の「[コンソールアクセスデバイス](#)」を参照してください。

Sun Cluster では専用の管理コンソールは必要ありませんが、専用の管理コンソールを使用すると、次のような利点があります。

- コンソールと管理ツールを同じマシンにまとめることで、クラスタ管理を一元化できます。
- ハードウェアサービスプロバイダによる問題解決が迅速に行われます。

管理コンソールの FAQ については、[第 4 章「よくある質問」](#)を参照してください。

## SPARC: Sun Cluster トポロジ

トポロジとは、Sun Cluster 環境で使用されている記憶装置プラットフォームにクラスタ内の Solaris ホストを接続するための接続スキームをいいます。Sun Cluster ソフトウェアは、次のガイドラインに従うトポロジをサポートします。

- SPARC ベースのシステムで構成される Sun Cluster 環境は、1つのクラスタで1つから16までの Solaris ホストをサポートします。ハードウェア構成によっては、SPARC ベースのシステムから成るクラスタで構成できるホストの最大数に制限が追加されます。
- 共有ストレージデバイスは、そのストレージデバイスでサポートされている数のホストに接続できます。
- 共有ストレージデバイスはクラスタのすべてのホストに接続する必要はありませんが、2つ以上のホストに接続する必要があります。

Logical Domains (LDoms) ゲストドメインおよび LDoms I/O ドメインを仮想 Solaris ホストとして構成できます。つまり、物理マシン、LDoms I/O ドメイン、および LDoms ゲストドメインの任意の組み合わせで構成される、クラスタペア、ペア +N 構成、N+1 構成、および N\*N 構成のクラスタを作成できます。また、LDoms ゲストドメインのみ、LDoms I/O ドメインのみ、またはこれら2つの組み合わせで構成されるクラスタを作成することもできます。

Sun Cluster ソフトウェアでは、特定のトポロジを使用するようにクラスタを構成する必要はありません。次のトポロジには、クラスタの接続スキームを説明するときを使用する用語を示します。これらのトポロジは典型的な接続スキームです。

- クラスタペア
- ペア +N
- N+1 (星型)
- N\*N (スケラブル)
- LDoms ゲストドメイン: ボックス内クラスタ
- LDoms ゲストドメイン: 2つのホストにわたる単一クラスタ
- LDoms ゲストドメイン: 2つのホストにわたる複数のクラスタ
- LDoms ゲストドメイン: 冗長 I/O ドメイン

次の各項では、それぞれのトポロジを図で示しています。

## SPARC: クラスタペアトポロジ

クラスタペアトポロジとは、単一のクラスタ管理フレームワークのもとで動作する複数の Solaris ホストペアをいいます。この構成では、ペアの間でのみフェイルオーバーが発生します。ただし、すべてのホストがクラスタインターコネクトによって接続されていて、Sun Cluster ソフトウェア制御のもとで動作します。このトポロジ

を使用する場合、1つのペアでパラレルデータベースアプリケーションを実行し、別のペアでフェイルオーバーまたはスケラブルなアプリケーションを実行できます。

クラスタファイルシステムを使用すると、2ペア構成も可能になります。アプリケーションデータが格納されているディスクにすべてのホストが直接接続されていない場合でも、複数のホストがスケラブルサービスまたはパラレルデータベースを実行できます。

次の図は、クラスタペア構成を示したものです。

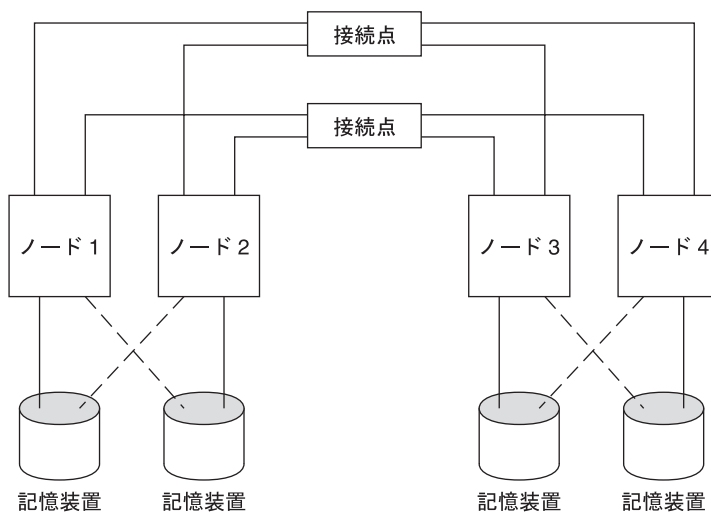


図2-2 SPARC: クラスタペアトポロジ

## SPARC: ペア +N トポロジ

ペア +N トポロジには、次のものに直接接続された Solaris ホストのペアが含まれています。

- 共有ストレージ。
- 共有ストレージにアクセスするため、クラスタインターコネクトを使用するホストの追加セット (それ自身は直接接続を持たない)。

次の図は、4つのホストのうち2つ (ホスト3とホスト4) がクラスタインターコネクトを使用して記憶装置にアクセスする、1つのペア +N トポロジを示したものです。この構成を拡張し、共有記憶装置には直接アクセスしない追加ホストを追加することができます。



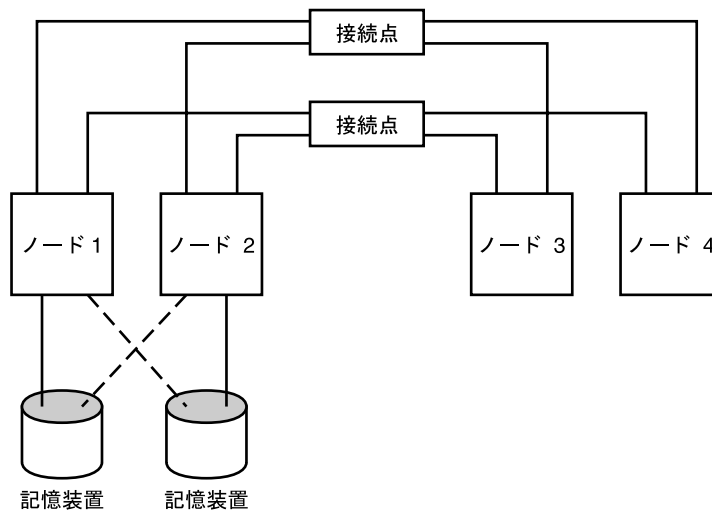


図 2-3 ペア +N トポロジ

## SPARC: N+1 (星形) トポロジ

N+1 トポロジには、複数の主 Solaris ホストと 1 つの二次ホストが含まれます。主ホストと二次ホストを同等に構成する必要はありません。主ホストは、アプリケーションサービスをアクティブに提供します。二次ホストは、主ホストに障害が生じるのを待機する間、アイドル状態である必要はありません。

二次ホストは、この構成ですべての多重ホスト記憶装置に物理的に接続されている唯一のホストです。

主ホストで障害が発生すると、Sun Cluster はそのリソースの処理を二次ホストで続行します。リソースは自動または手動で主ホストに切り換えられるまで二次ホストで機能します。

二次ホストには、主ホストの 1 つに障害が発生した場合に負荷を処理できるだけの十分な予備の CPU 容量が常に必要です。

次の図は、N+1 構成を示したものです。

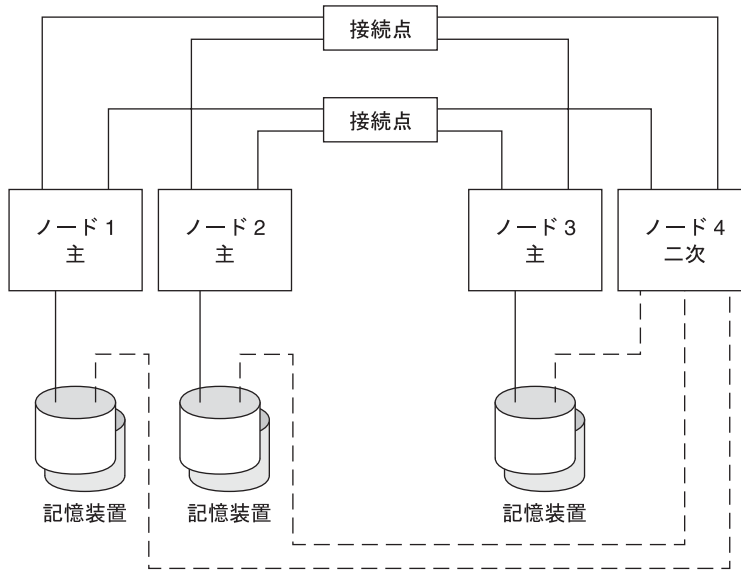


図 2-4 SPARC: N+1 トポロジ

## SPARC: N\*N (スケーラブル) トポロジ

N\*N トポロジを使用すると、クラスタ内のすべての共有ストレージデバイスをクラスタ内のすべての Solaris ホストに接続できます。このトポロジを使用すると、高可用性アプリケーションはサービスを低下させずに、あるホストから別のホストにフェイルオーバーできます。フェイルオーバーが発生すると、新しいホストはプライベートインターコネクトではなく、ローカルパスを使用して、記憶装置にアクセスできます。

次の図に、N\*N 構成を示します。

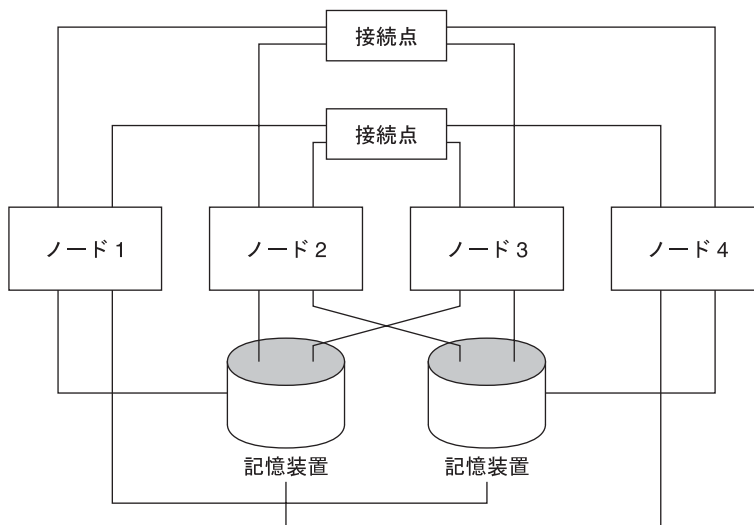


図 2-5 SPARC: N\*N トポロジ

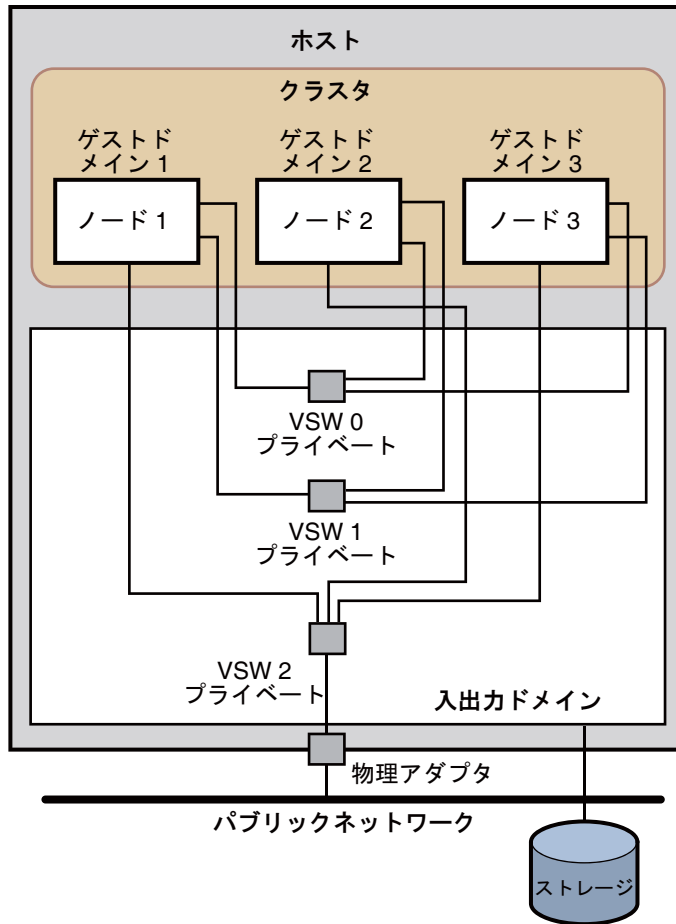
## SPARC: LDoms ゲストドメイン: ボックス内クラスタトポロジ

この Logical Domains (LDoms) ゲストドメイントポロジでは、クラスタと、そのクラスタ内のすべてのノードは同じ Solaris ホスト上に存在します。各 LDom ゲストドメインノードは、クラスタ内の Solaris ホストと同じように動作します。定足数デバイスを含める必要をなくすために、この構成には2つのノードのみではなく、3つのノードが含まれます。

このトポロジでは、プライベートネットワークの仮想スイッチ (vsw) 間の通信のみが必要であるため、各仮想スイッチを物理ネットワークに接続する必要はありません。このトポロジでは、すべてのクラスタノードが同じホスト上に置かれるため、クラスタノードは同じストレージデバイスを共有することもできます。クラスタ内の LDoms ゲストドメインおよび LDom I/O ドメインの使用法とインストールのガイドラインに関する詳細は、『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』の「Sun Logical Domains ソフトウェアをインストールしてドメインを作成する」を参照してください。

このトポロジでは、クラスタ内のすべてのノードが同じホスト上に置かれるため、高可用性は提供されません。ただし、開発者と管理者にとっては、テストや運用面以外の作業を行うのにこのトポロジが役立つ場合があります。このトポロジは、「ボックス内クラスタ」とも呼ばれます。

次の図はボックス内クラスタ構成を示したものです。



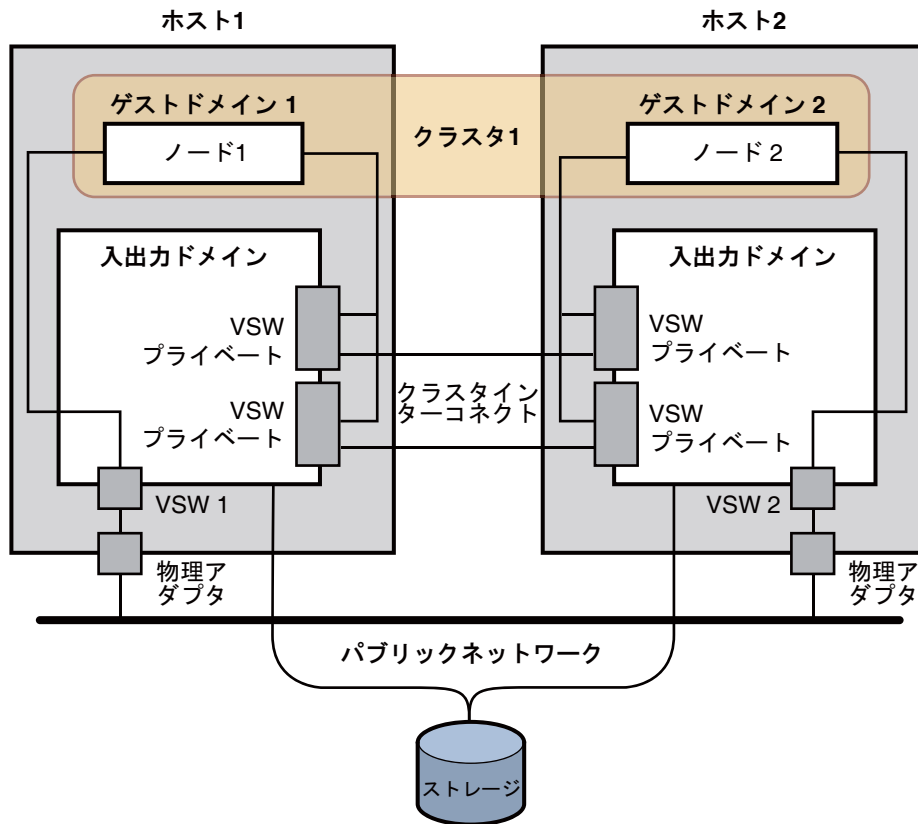
VSW = 仮想的な切り替え

図 2-6 SPARC: ボックス内クラスタトポロジ

## SPARC: LDoms ゲストドメイン: 2つのホストにわたる単一クラスタトポロジ

この Logical Domains (LDoms) ゲストドメイントポロジでは、単一クラスタは2つの異なる Solaris ホストにわたり、それぞれのクラスタは各ホスト上で1つのノードを構成します。各 LDoms ゲストドメインノードは、クラスタ内の Solaris ホストと同じように動作します。クラスタ内の LDoms ゲストドメインおよび LDoms I/O ドメインの使用法とインストールのガイドラインに関する詳細は、『[Sun Cluster ソフトウェアのインストール \(Solaris OS 版\)](#)』の「[Sun Logical Domains ソフトウェアをインストールしてドメインを作成する](#)」を参照してください。

次の図は、2つのホストにわたる単一クラスタの構成を示したものです。



VSW = 仮想的な切り替え

図 2-7 SPARC: 2つのホストにわたる単一クラスタ

## SPARC: LDoms ゲストドメイン: 2つのホストにわたる複数のクラスタトポロジ

この Logical Domains (LDoms) ゲストドメイントポロジでは、各クラスタは2つの異なる Solaris ホストにわたり、それぞれのクラスタは各ホスト上で1つのノードを構成します。各 LDom ゲストドメインノードは、クラスタ内の Solaris ホストと同じように動作します。この構成では、両方のクラスタで同じインターコネクトスイッチが共有されるため、各クラスタ上で別のプライベートネットワークアドレスを指定する必要があります。プライベートネットワークアドレスを変えずに、インターコネクトスイッチを共有するクラスタ上で同じプライベートネットワークアドレスを指定すると、構成に障害が発生します。

---

クラスタ内の LDoms ゲストドメインおよび LDoms I/O ドメインの使用法とインストールのガイドラインに関する詳細は、『[Sun Cluster ソフトウェアのインストール \(Solaris OS 版\)](#)』の「[Sun Logical Domains ソフトウェアをインストールしてドメインを作成する](#)」を参照してください。

次の図は、2つのホストにわたる2つ以上のクラスタの構成を示したものです。

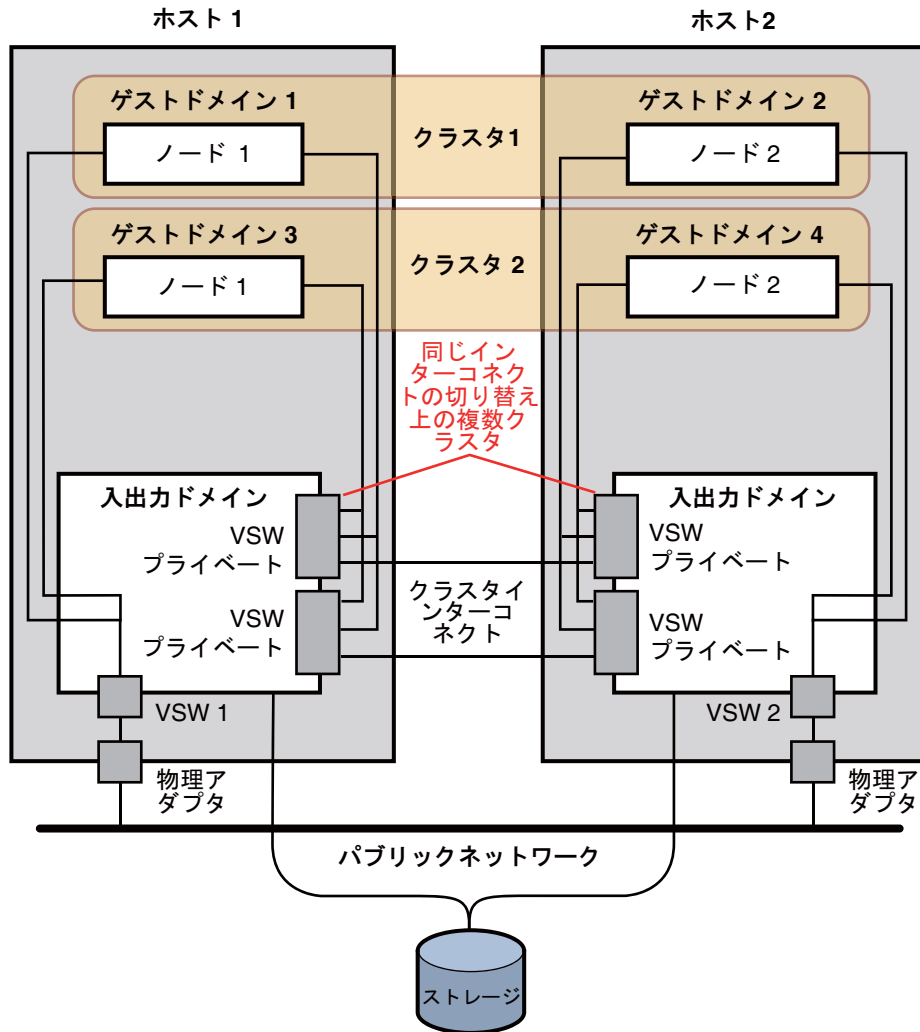


図 2-8 SPARC: 2つのホストにわたる複数のクラスタ

## SPARC: LDoms ゲストドメイン: 冗長 I/O ドメイン

この Logical Domains (LDoms) ゲストドメイントポロジでは、I/O ドメインで障害が発生した場合に、複数の I/O ドメインにより、ゲストドメインまたはクラスタ内のノードが動作し続けることを保証します。各 LDom ゲストドメインノードは、クラスタ内の Solaris ホストと同じように動作します。



このトポロジでは、ゲストドメインは、2つのパブリックネットワーク(1つは各I/Oドメインから)にわたってIPネットワークマルチパス化(IP network multipathing, IPMP)を実行します。また、ゲストドメインは、異なるI/Oドメインにわたってストレージデバイスをミラー化します。クラスタ内のLDomsゲストドメインおよびLDoms I/Oドメインの使用法とインストールのガイドラインに関する詳細は、『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』の「Sun Logical Domains ソフトウェアをインストールしてドメインを作成する」を参照してください。

次の図は、I/Oドメインで障害が発生した場合に、冗長I/Oドメインにより、クラスタ内のノードが動作し続けることを保証する構成を示したものです。

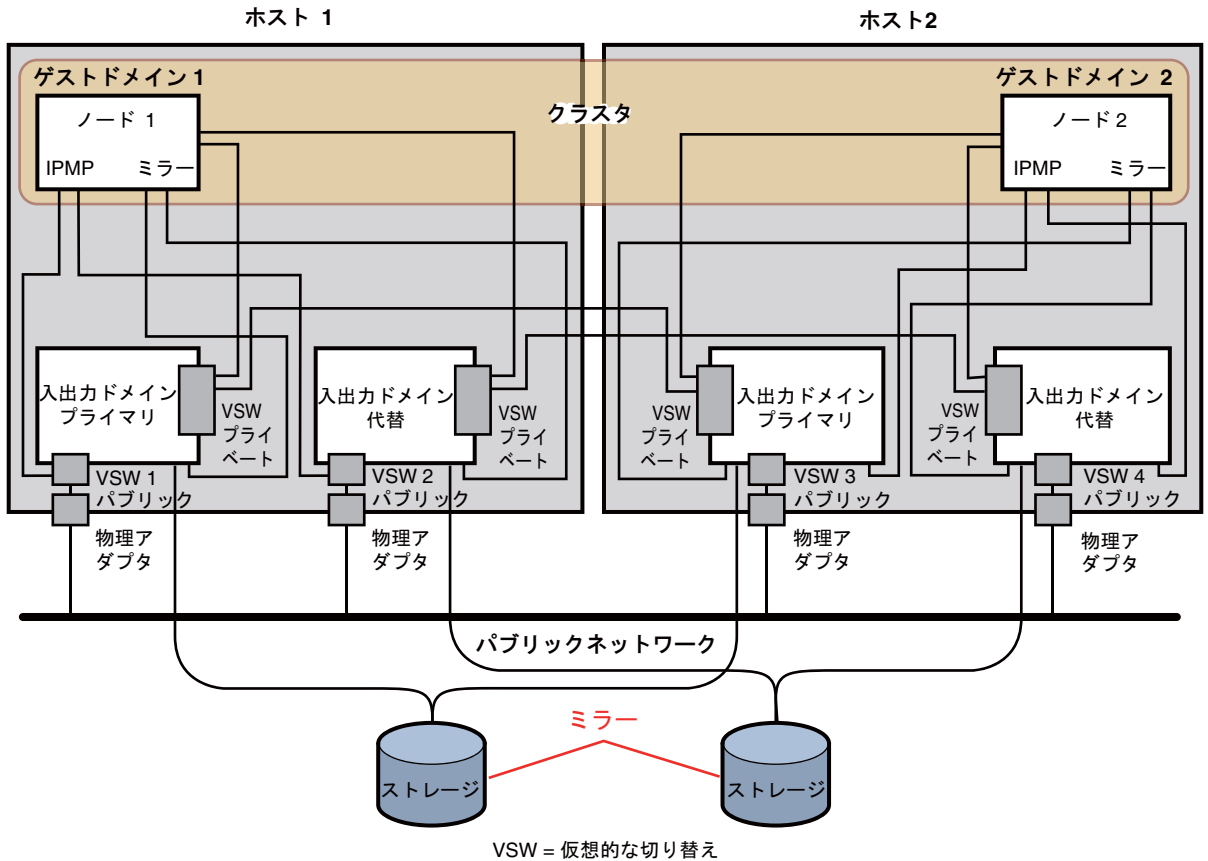


図 2-9 SPARC: 冗長 I/O ドメイン

## x86: Sun Cluster トポロジ

トポロジとは、クラスタノードと、クラスタで使用される記憶装置プラットフォームを接続する接続スキームをいいます。Sun Cluster は、次のガイドラインに従うトポロジをサポートします。

- Sun Cluster ソフトウェアは、1つのクラスタで1つから8つまでの Solaris ホストをサポートします。ハードウェア構成によっては、x86 ベースのシステムから成るクラスタで構成できるホストの最大数に制限が追加されます。サポートされるホスト構成については、42 ページの「x86: Sun Cluster トポロジ」を参照してください。
- 共有記憶装置をホストに接続する必要があります。

Sun Cluster では、特定のトポロジを使用するようにクラスタを構成する必要はありません。次のクラスタペアトポロジは、x86 ベースのホストから成るクラスタで可能なトポロジです。このトポロジを示すことによって、クラスタの接続スキームを表す用語を紹介します。このトポロジは代表的な接続スキームです。

次の項では、トポロジを図で示しています。

### x86: クラスタペアトポロジ

クラスタペアトポロジとは、単一のクラスタ管理フレームワークのもとで動作する2つの Solaris ホストをいいます。この構成では、ペアの間でのみフェイルオーバーが発生します。ただし、すべてのホストがクラスタインターコネクトによって接続されていて、Sun Cluster ソフトウェア制御のもとで動作します。このトポロジを使用する場合、ペアでパラレルデータベース、フェイルオーバー、またはスケラブルアプリケーションを実行できます。

次の図は、クラスタペア構成を示したものです。

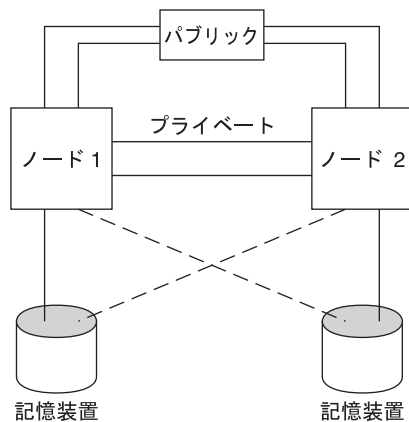


図 2-10 x86: クラスタペアドトポロジ

## x86: N+1 (星形) トポロジ

N+1 トポロジには、複数の主 Solaris ホストと 1 つの二次ホストが含まれます。主ホストと二次ホストを同等に構成する必要はありません。主ホストは、アプリケーションサービスをアクティブに提供します。二次ホストは、主ホストに障害が生じるのを待機する間、アイドル状態である必要はありません。

二次ホストは、この構成ですべての多重ホスト記憶装置に物理的に接続されている唯一のホストです。

主ホストで障害が発生すると、Sun Cluster はそのリソースの処理を二次ホストで続行します。リソースは自動または手動で主ホストに切り換えられるまで二次ホストで機能します。

二次ホストには、主ホストの 1 つに障害が発生した場合に負荷を処理できるだけの十分な予備の CPU 容量が常に必要です。

次の図は、N+1 構成を示したものです。

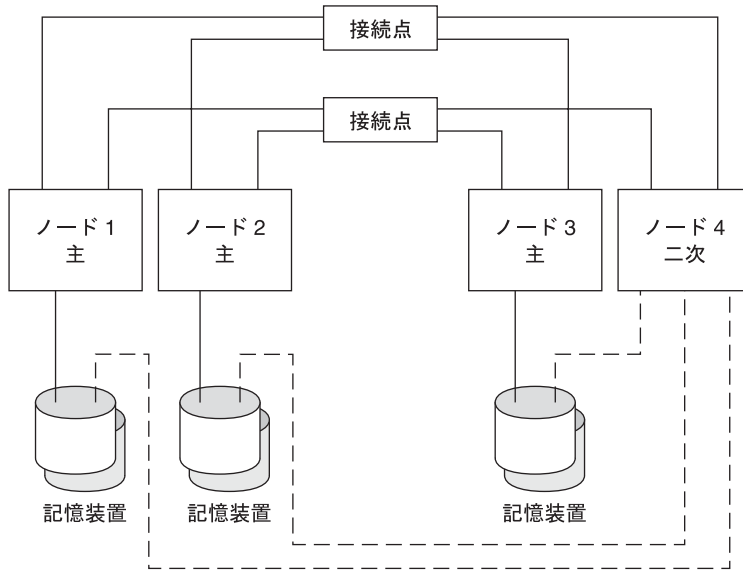


図 2-11 x86: N+1 トポロジ

## 重要な概念 - システム管理者とアプリケーション開発者

---

この章では、Sun Cluster 環境のソフトウェアコンポーネントに関する重要な概念について説明します。この章の情報は、主に Sun Cluster API および SDK を使用するシステム管理者およびアプリケーション開発者向けです。クラスタの管理者にとっては、この情報は、クラスタソフトウェアのインストール、構成、管理についての予備知識となります。アプリケーション開発者は、この情報を使用して、作業を行うクラスタ環境を理解できます。

この章の内容は次のとおりです。

- 46 ページの「管理インタフェース」
- 46 ページの「クラスタ内の時間」
- 47 ページの「高可用性フレームワーク」
- 50 ページの「グローバルデバイス」
- 51 ページの「デバイスグループ」
- 55 ページの「広域名前空間」
- 56 ページの「クラスタファイルシステム」
- 60 ページの「ディスクパスの監視」
- 63 ページの「定足数と定足数デバイス」
- 71 ページの「データサービス」
- 80 ページの「新しいデータサービスの開発」
- 82 ページの「クラスタインターコネクトによるデータサービストラフィックの送受信」
- 84 ページの「リソース、リソースグループ、リソースタイプ」
- 88 ページの「Solaris ゾーンをサポート」
- 91 ページの「サービス管理機能」
- 92 ページの「システムリソースの使用状況」
- 95 ページの「データサービスプロジェクトの構成」
- 105 ページの「パブリックネットワークアダプタと IP ネットワークマルチパス」
- 106 ページの「SPARC: 動的再構成のサポート」

## 管理インタフェース

複数のユーザーインタフェースから Sun Cluster ソフトウェアをインストール、構成、および管理する方法を選択できます。Sun Cluster Manager のグラフィカルユーザーインタフェース (GUI) またはコマンドラインインタフェースのいずれかによって、システム管理作業を実行できます。コマンド行インタフェースでは、特定のインストール作業や構成作業を容易にする `scinstall` や `clsetup` などのユーティリティーが使用できます。Sun Cluster ソフトウェアには、Sun Management Center の一部として実行される、特定のクラスタ作業に GUI を提供するモジュールもあります。このモジュールを使用できるのは、SPARC ベースのクラスタに限られます。管理インタフェースについての詳細は、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の「管理ツール」を参照してください。

## クラスタ内の時間

クラスタ内のすべての Solaris ホスト間の時刻は同期をとる必要があります。クラスタホストの時刻と外部の時刻ソースの同期をとるかどうかは、クラスタの操作にとって重要ではありません。Sun Cluster ソフトウェアは、時間情報プロトコル (Network Time Protocol、NTP) を使用し、ホスト間のクロックの同期をとっています。

通常、システムクロックが数分の 1 秒程度変更されても問題は起こりません。しかし、システムクロックと時刻の起点の同期をとるために、`date`、`rdate`、`xntpdate` を (対話形式または `cron` スクリプト内で) アクティブクラスタに対して実行すると、これよりも大幅な時刻変更を強制的に行うことが可能です。ただしこの強制的な変更を行った場合、ファイル修正時刻の表示に問題が生じたり、NTP サービスに混乱が生じる可能性があります。

Solaris オペレーティングシステムを各クラスタホストにインストールする場合は、ホストのデフォルトの時刻と日付の設定を変更できます。通常は、工場出荷時のデフォルト値を使用します。

`scinstall` コマンドを使用して Sun Cluster ソフトウェアをインストールする場合は、インストールプロセスの手順の 1 つとして、クラスタの NTP を構成します。Sun Cluster ソフトウェアは、`ntp.cluster` というテンプレートファイルを提供しています (インストールされたクラスタホストの `/etc/inet/ntp.cluster` を参照)。このテンプレートは、すべてのクラスタホスト間で対等関係を確立します。1 つのホストは「優先ホスト」になります。ホストはプライベートホスト名で識別され、時刻の同期化がクラスタインターコネクト全体で行われます。NTP 用のクラスタの構成方法については、『[Sun Cluster ソフトウェアのインストール \(Solaris OS 版\)](#)』の第 2 章「グローバルクラスタノードへのソフトウェアのインストール」を参照してください。

また、クラスタの外部に 1 つまたは複数の NTP サーバーを設定し、`ntp.conf` ファイルを変更してその構成を反映させることもできます。

通常の操作では、クラスタの時刻を調整する必要はありません。ただし、Solaris オペレーティングシステムをインストールしたときに設定された誤った時刻を変更する場合の手順については、『Sun Cluster のシステム管理 (Solaris OS 版)』の第 8 章「クラスタの管理」を参照してください。

## 高可用性フレームワーク

Sun Cluster ソフトウェアでは、ユーザーとデータ間の「パス」にあるすべてのコンポーネント、つまり、ネットワークインタフェース、アプリケーション自体、ファイルシステム、および多重ホストデバイスを高可用性にします。一般に、システムで単一(ソフトウェアまたはハードウェア)の障害が発生してもあるクラスタコンポーネントが稼働し続けられる場合、そのコンポーネントは高可用性であると考えられます。

次の表は Sun Cluster コンポーネント障害の種類(ハードウェアとソフトウェアの両方)と高可用性フレームワークに組み込まれた回復の種類を示しています。

表 3-1 Sun Cluster の障害の検出と回復のレベル

| 障害が発生したクラス<br>タリソース | ソフトウェアの回復   | ハードウェアの回復                              |
|---------------------|---|--|
| データサービス             | HA API、HA フレームワーク   | なし                                     |
| パブリックネットワークアダプタ     | IP ネットワークマルチパス  | 複数のパブリックネットワークアダプタカード                  |
| クラスタファイルシステム        | 一次複製と二次複製   | 多重ホストデバイス                              |
| ミラー化された多重ホストデバイス    | ボリューム管理 (Solaris Volume Manager および Veritas Volume Manager) | ハードウェア RAID-5 (Sun StorEdge™ A3x00 など) |
| 広域デバイス              | 一次複製と二次複製   | デバイス、クラスタトランスポート接続点への多重パス              |
| プライベートネットワーク        | HA トランスポートソフトウェア  | ハードウェアから独立した多重プライベートネットワーク             |
| ホスト                 | CMM、フェイルファーストドライバ   | 複数ホスト                                  |
| ゾーン                 | HA API、HA フレームワーク   | なし                                     |

Sun Cluster ソフトウェアの高可用性フレームワークは、ノードの障害をすばやく検出して、クラスタ内の残りのノードにあるフレームワークリソース用に新しい同等のサーバーを作成します。どの時点でもすべてのフレームワークリソースが使用できなくなることはありません。障害が発生したノードの影響を受けないフレームワ

ークリソースは、回復中も完全に使用できます。さらに、障害が発生したノードのフレームワークリソースは、回復されると同時に使用可能になります。回復されたフレームワークリソースは、ほかのすべてのフレームワークリソースが回復するまで待機する必要はありません。

最も可用性の高いフレームワークリソースは、そのリソースを使用するアプリケーション(データサービス)に対して透過的に回復されます。フレームワークリソースのアクセス方式は、ノードの障害時にも完全に維持されます。アプリケーションは、フレームワークリソースサーバーが別のノードに移動したことを認識できないだけです。1つのノードで障害が発生しても、残りのノード上にあるプログラムがそのノードのファイル、デバイス、およびディスクボリュームを使用できるので、その障害は完全に透過的と言えます。別のホストからそのディスクに代替ハードウェアパスが設定されている場合に、このような透過性が実現されます。この例としては、複数ホストへのポートを持つ多重ホストデバイスの使用があります。

## ゾーンメンバーシップ

また、Sun Cluster ソフトウェアはゾーンがいつ起動または停止するかを検出することによって、ゾーンメンバーシップを追跡します。こうした変化も再構成の原因となります。再構成によって、クラスタ内のノード間にクラスタリソースを再配置できます。

## クラスタメンバーシップモニター

データが破壊から保護されるように保証するには、すべてのノードが、クラスタメンバーシップに対して一定の同意に達していなければなりません。必要であれば、CMMは、障害に応じてクラスタサービス(アプリケーション)のクラスタ再構成を調整します。

CMMは、クラスタのトランスポート層から、他のノードへの接続に関する情報を受け取ります。CMMは、クラスタインターコネクトを使用して、再構成中に状態情報を交換します。

CMMは、クラスタメンバーシップの変更を検出すると、それに合わせてクラスタを構成します。このような同期構成では、クラスタの新しいメンバーシップに基づいて、クラスタリソースが再配布されることがあります。

## フェイルファースト機構

「フェイルファースト」機構では、グローバルクラスタ投票ノードまたはグローバルクラスタ非投票ノードのいずれかにおける重大な問題が検出されます。フェイルファーストで問題が検出されたときに、Sun Cluster が取る措置は、問題が投票ノードで発生するか非投票ノードで発生するかによって異なります。



重大な問題が投票ノードで発生した場合、Sun Cluster は強制的にノードを停止させます。Sun Cluster は次にノードをクラスタメンバーシップから削除します。

重大な問題が非投票ノードで発生した場合、Sun Cluster は非投票ノードを再起動します。

ノードは、ほかのノードとの接続を失うと、通信が可能なノードとクラスタを形成しようとしています。そのセットのノードが定足数に達しない場合、Sun Cluster ソフトウェアはノードを停止して、共有ディスクからノードをフェンスします。つまり、ノードの共有ディスクへのアクセスを遮ります。

フェンシングは、選択したディスクまたはすべてのディスクに対してオフにできません。



注意 - 不適切な状況でフェンシングを無効にすると、アプリケーションのフェイルオーバー時にデータが破損する危険性が高くなります。フェンシングの無効化を検討する場合には、データ破損の可能性を十分に調査してください。SATA (Serial Advanced Technology Attachment) ディスクなど、共有記憶装置が SCSI プロトコルに対応していない場合、またはクラスタの外部にあるホストからクラスタの記憶装置へのアクセスを許可する場合にフェンシングをオフにします。

1つまたは複数のクラスタ固有のデーモンが停止すると、Sun Cluster ソフトウェアは重大な問題が発生したことを宣言します。Sun Cluster ソフトウェアは、投票ノードと非投票ノードの両方でクラスタ固有のデーモンを実行します。重大な問題が発生すると、Sun Cluster はノードを停止して削除するか、問題が発生した非投票ノードを再起動します。

非投票ノードで実行されるクラスタ固有のデーモンが失敗すると、次のようなメッセージがコンソールに表示されます。

```
cl_runtime: NOTICE: Failfast: Aborting because "pmfd" died in zone "zone4" (zone id 3)
35 seconds ago.
```

投票ノードで実行されるクラスタ固有のデーモンが失敗し、ノードでパニックが発生すると、次のようなメッセージがコンソールに表示されます。

```
panic[cpu1]/thread=2a10007fcc0: Failfast: Aborting because "pmfd" died in zone "global" (zone id 0)
35 seconds ago.
409b8 cl_runtime: _0FZsc_syslog_msg_log_no_argsPviTCPCcTB+48 (70f900, 30, 70df54, 407acc, 0)
%l0-7: 1006c80 000000a 000000a 10093bc 406d3c80 7110340 0000000 4001 fbfd
```

パニック後、Solaris ホストは再起動し、ノードはクラスタに再び参加しようとする場合があります。あるいは、SPARC ベースのシステムで構成されているクラスタの場合、そのホストは OpenBoot PROM (OBP) プロンプトのままになる場合があります。ホストがどちらのアクションをとるかは、`auto-boot?` パラメータの設定によって

決定されます。OpenBoot PROM の ok プロンプトで、eeprom コマンドにより auto-boot? を設定できます。詳細は、[eeprom\(1M\)](#) のマニュアルページを参照してください。

## クラスタ構成レポジトリ (CCR)

CCR は、更新に 2 フェーズのコミットアルゴリズムを使用します。更新はすべてのクラスタメンバーで正常に終了する必要があります。そうしないと、その更新はロールバックされます。CCR はクラスタインターコネクトを使用して、分散更新を適用します。



注意 - CCR はテキストファイルで構成されていますが、CCR ファイルは絶対に自分では編集しないでください。各ファイルには、ノード間の一貫性を保証するための検査合計レコードが含まれています。CCR ファイルを自分で更新すると、ノードまたはクラスタ全体の機能が停止する可能性があります。

CCR は、CMM に依存して、定足数 (quorum) が確立された場合にのみクラスタが実行されるように保証します。CCR は、クラスタ全体のデータの一貫性を確認し、必要に応じて回復を実行し、データへの更新を容易にします。

## グローバルデバイス

Sun Cluster ソフトウェアは、「広域デバイス」を使用して、デバイスが物理的に接続されている場所に関係なく、任意のノードからクラスタ内のすべてのデバイスに対して、クラスタ全体の可用性の高いアクセスを可能にします。通常、広域デバイスへのアクセスを提供中のノードにエラーが発生すると、Sun Cluster ソフトウェアは自動的にこのデバイスへの別のパスを見つけます。Sun Cluster ソフトウェアは、アクセスをこのパスにリダイレクトさせます。Sun Cluster 広域デバイスには、ディスク、CD-ROM、テープが含まれます。しかし、Sun Cluster ソフトウェアがサポートする多重ポート広域デバイスはディスクだけです。つまり、CD-ROM とテープは現在、高可用性のデバイスではありません。各サーバーのローカルディスクも多重ポート化されていないため、可用性の高いデバイスではありません。

クラスタは、クラスタ内の各ディスク、CD-ROM、テープデバイスに一意の ID を自動的に割り当てます。この割り当てによって、クラスタ内の任意のノードから各デバイスに対して一貫したアクセスが可能になります。広域デバイス名前空間は、`/dev/global` ディレクトリにあります。詳細は、[55 ページの「広域名前空間」](#)を参照してください。

多重ポート広域デバイスは、1つのデバイスに対して複数のパスを提供します。多重ホストディスクは複数の Solaris ホストによってホストされるデバイスグループの一部であるため、多重ホストディスクの可用性は高くなります。

## デバイス ID と DID 疑似ドライバ

Sun Cluster ソフトウェアは、DID 疑似ドライバと呼ばれる構造によって広域デバイスを管理します。このドライバを使用して、多重ホストディスク、テープドライブ、CD-ROM を含め、クラスタ内のあらゆるデバイスに一意の ID を自動的に割り当てます。

DID 疑似ドライバは、クラスタの広域デバイスアクセス機能における重要な部分です。DID ドライバは、クラスタのすべてのノードを探索して、一意のデバイスのリストを作成し、クラスタのすべてのノードで一貫している一意のメジャー番号およびマイナー番号を各デバイスに割り当てます。広域デバイスへのアクセスは、ディスクを示す `c0t0d0` などの従来の Solaris デバイス ID ではなく、(DID ドライバが割り当てた)この一意のデバイス ID を利用して行われます。

この方法により、ディスクにアクセスするすべてのアプリケーション (ボリュームマネージャーまたは raw デバイスを使用するアプリケーションなど) は、一貫したパスを使用してクラスタ全体にアクセスできます。各デバイスのローカルメジャー番号およびマイナー番号は Solaris ホストによって異なり、Solaris デバイス命名規則も変更する可能性があるため、この一貫性は、多重ホストディスクにとって特に重要です。たとえば、Host1 は多重ホストディスクを `c1t2d0` と識別し、Host2 は同じディスクをまったく異なるディスクとして、つまり、`c3t2d0` と識別する場合があります。ホストはこのような名前の代わりに、DID ドライバが割り当てた広域名 (`d10` など) を使用します。つまり、DID ドライバは多重ホストディスクへの一貫したマッピングを各ホストに提供します。

`cldevice` コマンドでデバイス ID を更新して管理します。詳細は、[cldevice\(1CL\)](#) のマニュアルページを参照してください。

## デバイスグループ

Sun Cluster ソフトウェアでは、多重ホストデバイスをすべて Sun Cluster ソフトウェアで管理する必要があります。最初に多重ホストディスク上にボリュームマネージャーのディスクグループ (Solaris Volume Manager のディスクセットまたは Veritas Volume Manager ディスクグループのいずれか) を作成します。次に、ボリュームマネージャーのディスクグループを「デバイスグループ」として登録します。デバイスグループは、広域デバイスの一種です。さらに、Sun Cluster ソフトウェアは、個々のディスクデバイスやテープデバイスごとに raw デバイスグループを自動的に作成します。ただし、これらのクラスタデバイスグループは、広域デバイスとしてアクセスされるまではオフラインの状態になっています。

この登録によって、Sun Cluster ソフトウェアは、どの Solaris ホストがどのボリュームマネージャーディスクグループへのパスをもっているかを知ることができます。この時点でそのボリュームマネージャーデバイスグループは、クラスタ内で広域アクセスが可能になります。あるデバイスグループが複数のホストから書き込み可能

(マスター)な場合は、そのデバイスグループに格納されるデータは、高度な可用性を有することになります。高度な可用性を備えたデバイスグループには、クラスタファイルシステムを格納できます。

---

注-デバイスグループは、リソースグループとは別のものです。あるノードは、データサービスプロセスのグループを表すリソースグループをマスターすることができます。別のノードは、データサービスによりアクセスされているディスクグループをマスターすることができます。ただし、もっとも良い方法は、特定のアプリケーションのデータを保存するデバイスグループと、アプリケーションのリソース(アプリケーションデーモン)を同じノードに含むリソースグループを維持することです。デバイスグループとリソースグループの関係についての詳細は、『[Sun Cluster データサービスの計画と管理 \(Solaris OS 版\)](#)』の「[リソースグループとデバイスグループの関係](#)」を参照してください。

---

あるノードがディスクデバイスグループを使用するとき、ボリュームマネージャーのディスクグループは実際に使用するディスクに対してマルチパスサポートを提供するため、そのディスクグループは「広域」になります。多重ホストディスクに物理的に接続された各クラスタホストは、デバイスグループへのパスを提供します。

## デバイスグループのフェイルオーバー

ディスク格納装置は複数の Solaris ホストに接続されるため、現在デバイスグループをマスターしているホストに障害が生じた場合でも、代替パスによってその格納装置にあるすべてのデバイスグループにアクセスできます。デバイスグループをマスターするホストの障害は、回復と一貫性の検査を実行するために要する時間を除けば、デバイスグループへのアクセスに影響しません。この時間の間は、デバイスグループが使用可能になるまで、すべての要求は(アプリケーションには透過的に)阻止されます。

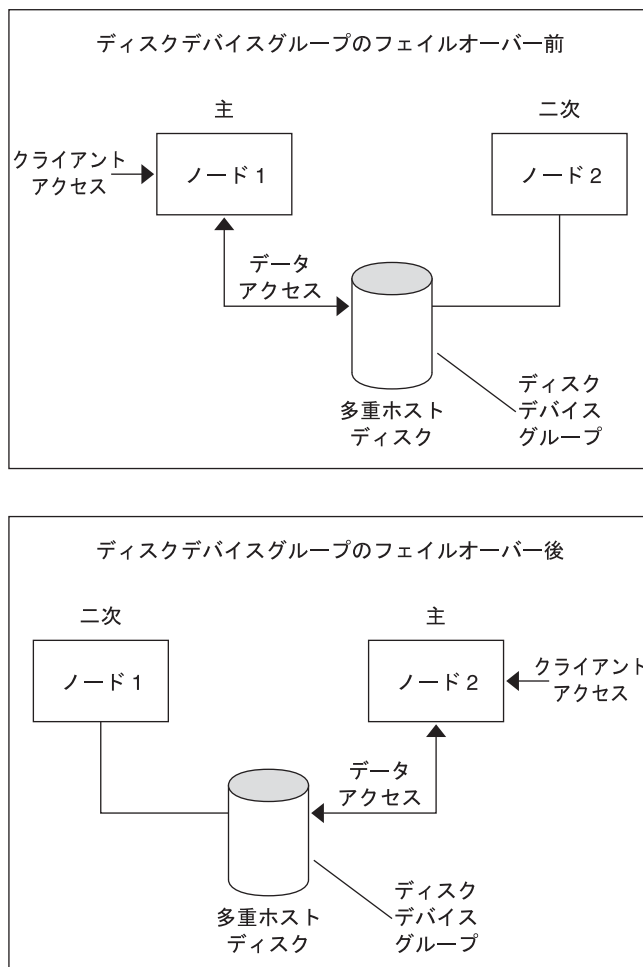


図 3-1 フェイルオーバー前後のデバイスグループ

## 多重ポートデバイスグループ

この節では、多重ポートディスク構成において性能と可用性をバランスよく実現するディスクデバイスグループのプロパティについて説明します。Sun Cluster ソフトウェアには、多重ポートディスク構成を設定するためのプロパティが2つあります。つまり、`preferenced` と `numsecondaries` です。`preferenced` プロパティは、フェイルオーバーの発生時に各ノードがどの順で制御を取得するかを制御します。`numsecondaries` プロパティは、特定のデバイスグループに対する二次ノードの数を設定します。

高可用性サービスは、主ノードが停止し、かつ、主ノードになる資格のある二次ノードが存在しないときに、完全に停止したと見なされます。preferenced プロパティが true に設定されている場合、サービスのフェイルオーバーが発生すると、ノードリストの順序に従って二次ノードが選択されます。ノードリストは、ノードが主制御を引き受ける順序、またはスベアから二次への移行を引き受ける順序を決めます。clsetup コマンドを使用して、デバイスサービスの優先順序を動的に変更できます。従属サービスプロバイダ(広域ファイルシステムなど)に関連する設定は、デバイスサービスの設定と同じになります。

主ノードは、正常な運用時に二次ノードのチェックポイントをとります。多重ポートディスク構成では、二次ノードのチェックポイントをとるたびに、クラスタの性能の低下やメモリーのオーバーヘッドの増加が発生します。スベアノードのサポートが実装されているのは、このようなチェックポイントによる性能の低下やメモリーのオーバーヘッドを最小限に抑えるためです。デフォルトでは、デバイスグループには1つの主ノードと1つの二次ノードがあります。残りのプロバイダノードはスベアノードです。フェイルオーバーが発生すると、二次ノードが主ノードになり、ノードリスト上でもっとも優先順位の高い(スベア)ノードが二次ノードになります。

二次ノードの望ましい数には、任意の整数(1から、デバイスグループ内の動作可能な主ノード以外のプロバイダノードの数まで)を設定できます。

---

注 - Solaris Volume Manager を使用している場合、numsecondaries プロパティにデフォルト以外の数字を設定するには、まず、デバイスグループを作成する必要があります。

---

デバイスサービスの二次ノードのデフォルト数は1です。望ましい数とは、複製フレームワークによって維持される二次プロバイダノードの実際の数です。ただし、動作可能な主ノード以外のプロバイダノードの数が望ましい数よりも小さい場合を除きます。ノードを構成に追加したり、ノードを構成から切り離す場合は、numsecondaries プロパティを変更したあと、ノードリストを十分に確認する必要があります。ノードリストと二次ノードの数を正しく保つことによって、構成されている二次ノードの数と、フレームワークによって与えられている実際の数の不一致を防げます。

- (Solaris Volume Manager) 構成へのノードの追加および構成からのノードの削除を管理するには、Solaris Volume Manager デバイスグループ用の metaset コマンドを preferred および numsecondaries プロパティ設定と組み合わせて使用します。
- (Veritas Volume Manager) 構成へのノードの追加および構成からのノードの削除を管理するには、VxVM デバイスグループ用の cldevicegroup コマンドを preferred および numsecondaries プロパティ設定と組み合わせて使用します。

- デバイスグループのプロパティの変更手順については、『Sun Cluster のシステム管理 (Solaris OS 版)』の「クラスタファイルシステムの管理の概要」を参照してください。

## 広域名前空間

広域デバイスを有効にする Sun Cluster ソフトウェアの機構は、「広域名前空間」です。広域名前空間には、ボリューム管理ソフトウェアの名前空間とともに、`/dev/global/` 階層が含まれます。広域名前空間は、多重ホストディスクとローカルディスクの両方 (および CD-ROM やテープなどのほかのクラスタデバイスすべて) を反映して、多重ホストディスクへの複数のフェイルオーバーパスを提供します。多重ホストディスクに物理的に接続された各 Solaris ホストは、クラスタ内のすべてのノードの記憶装置に対するパスを提供します。

Solaris Volume Manager の場合、ボリュームマネージャーの名前空間は、通常、`/dev/md/diskset/dsk` (と `rdsk`) ディレクトリにあります。Veritas VxVM の場合、ボリュームマネージャーの名前空間は `/dev/vx/dsk/disk-group` ディレクトリと `/dev/vx/rdsk/disk-group` ディレクトリにあります。これらの名前空間は、クラスタ全体でインポートされている Solaris Volume Manager の各ディスクセットと VxVM の各ディスクグループのディレクトリから構成されます。これらの各ディレクトリには、そのディスクセットまたはディスクグループ内の各メタデバイスまたはボリュームのデバイスホストが格納されています。

Sun Cluster ソフトウェアでは、ローカルボリュームマネージャーの名前空間内の各デバイスホストは `/global/.devices/node@nodeID` ファイルシステム内のデバイスホストへのシンボリックリンクに置き換えられます。`nodeID` は、クラスタ内のノードを表す整数です。Sun Cluster ソフトウェアは、その標準的な場所に引き続きシンボリックリンクとしてボリューム管理デバイスも表示します。広域名前空間と標準ボリュームマネージャー名前空間は、どちらも任意のクラスタノードから使用できます。

広域名前空間には、次の利点があります。

- 各ホストの独立性が高く、デバイス管理モデルを変更する必要がほとんどありません。
- デバイスを選択的に広域に設定できます。
- Sun の製品以外のリンクジェネレータが引き続き動作します。
- ローカルデバイス名を指定すると、その広域名を取得するために簡単なマッピングが提供されます。

## ローカル名前空間と広域名前空間の例

次の表は、多重ホストディスク `c0t0d0s0` でのローカル名前空間と広域名前空間のマッピングを示したものです。

表 3-2 ローカル名前空間と広域名前空間のマッピング

| コンポーネントまたはパス           | ローカルホスト名前空間                            | 広域名前空間   |
|------------------------|--|--|
| Solaris 論理名            | <code>/dev/dsk/c0t0d0s0</code>         | <code>/global/.devices/node@nodeID /dev/dsk/c0t0d0s0</code>          |
| DID 名                  | <code>/dev/did/dsk/d0s0</code>         | <code>/global/.devices/node@nodeID /dev/did/dsk/d0s0</code>          |
| Solaris Volume Manager | <code>/dev/md/diskset/dsk/d0</code>    | <code>/global/.devices/node@nodeID /dev/md/diskset/dsk/d0</code>     |
| Veritas Volume Manager | <code>/dev/vx/dsk/disk-group/v0</code> | <code>/global/.devices/node@nodeID /dev/vx/dsk/disk-group /v0</code> |

広域名前空間はインストール時に自動的に生成されて、再構成再起動のたびに更新されます。広域名前空間は、`cldevice` コマンドを使用して生成することもできます。詳細は、[cldevice\(1CL\)](#) のマニュアルページを参照してください。

## クラスタファイルシステム

クラスタファイルシステムには、次の機能があります。

- ファイルのアクセス場所が透過的になります。システムのどこにあるファイルでも、プロセスから開くことができます。すべての Solaris ホストのプロセスから同じパス名を使ってファイルにアクセスできます。

---

注-クラスタファイルシステムは、ファイルを読み取る際に、ファイル上のアクセス時刻を更新しません。

---

- 一貫したプロトコルを使用して、ファイルが複数のノードから同時にアクセスされている場合でも、UNIX ファイルアクセス方式を維持します。
- 拡張キャッシュ機能とゼロコピーバルク入出力移動機能により、ファイルデータを効率的に移動することができます。
- クラスタファイルシステムには、`fcntl` コマンドインタフェースに基づく、高度な可用性を備えたアドバイザリファイルロッキング機能があります。クラスタファイルシステムに対してアドバイザリファイルロッキング機能を使えば、複数のクラスタノードで動作するアプリケーションの間で、データのアクセスを同期化で



きます。ファイルロックを所有するノードがクラスタから切り離されたり、ファイルロックを所有するアプリケーションが異常停止すると、それらのロックはただちに解放されます。

- 障害が発生した場合でも、データへの連続したアクセスが可能です。アプリケーションは、ディスクへのパスが有効であれば、障害による影響を受けません。この保証は、raw ディスクアクセスとすべてのファイルシステム操作で維持されます。
- クラスタファイルシステムは、基本のファイルシステムからもボリュームマネージャーからも独立しています。クラスタシステムファイルは、サポートされているディスク上のファイルシステムすべてを広域にします。

広域デバイスにファイルシステムをマウントするとき、広域にマウントする場合は `mount -g` を使用し、ローカルにマウントする場合は `mount` を使用します。

プログラムは、同じファイル名(たとえば、`/global/foo`)によって、クラスタ内のすべてのノードからクラスタファイルシステムのファイルにアクセスできます。

クラスタファイルシステムは、すべてのクラスタメンバーにマウントされます。クラスタファイルシステムをクラスタメンバーのサブセットにマウントすることはできません。

クラスタファイルシステムは、特定のファイルシステムタイプではありません。つまり、クライアントは、実際に使用するファイルシステム(UFSなど)だけを認識します。

## クラスタファイルシステムの使用法

Sun Cluster ソフトウェアでは、すべての多重ホストディスクがディスクデバイスグループとして構成されています。これは、Solaris Volume Manager のディスクセット、VxVM のディスクグループ、またはソフトウェアベースのボリューム管理ソフトウェアの制御下でない個々のディスクが該当します。

クラスタファイルシステムを高可用性にするには、使用するディスクストレージが複数の Solaris ホストに接続されていなければなりません。したがって、ローカルファイルシステム(ホストのローカルディスクに格納されているファイルシステム)をクラスタファイルシステムにした場合は、高可用性にはなりません。

クラスタファイルシステムは、通常のファイルシステムと同様にマウントできます。

- 手作業。mount コマンドと -g または -o global マウントオプションを使用し、コマンド行からクラスタファイルシステムをマウントします。次に例を示します。

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

- 自動。マウントオプションで /etc/vfstab ファイルにエントリを作成して、ブート時にクラスタファイルシステムをマウントします。さらに、すべてのホストの /global ディレクトリ下にマウントポイントを作成します。ディレクトリ /global を推奨しますが、ほかの場所でも構いません。次に、/etc/vfstab ファイルの、クラスタファイルシステムを示す行の例を示します。

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

---

注 - Sun Cluster ソフトウェアには、クラスタファイルシステムに対する特定の命名規則はありません。しかし、/global/disk-group などのように、同じディレクトリのもとにすべてのクラスタファイルシステムのマウントポイントを作成すると、管理が容易になります。詳細は、『Sun Cluster 3.1 9/04 Software Collection for Solaris OS (SPARC Platform Edition)』と『Sun Cluster のシステム管理 (Solaris OS 版)』を参照してください。

---

## HASStoragePlus リソースタイプ

HASStoragePlus リソースタイプは、ローカルおよび広域ファイルシステム構成を高可用対応にするように設計されています。HASStoragePlus リソースタイプを使用して、ローカルまたは広域ファイルシステムを Sun Cluster 環境に統合し、このファイルシステムを高可用対応にすることができます。

HASStoragePlus リソースタイプを使用すると、ファイルシステムをグローバルクラスタ非投票ノードで利用可能にすることができます。HASStoragePlus リソースタイプを使用してこのようにするには、グローバルクラスタ投票ノードとグローバルクラスタ非投票ノードにマウントポイントを作成してください。ファイルシステムをグローバルクラスタ非投票ノードで利用可能にするために、HASStoragePlus リソースタイプは、まずグローバルクラスタ投票ノードにあるファイルシステムをマウントします。このリソースタイプは、次にグローバルクラスタ非投票ノードでループバックマウントを実行します。

---

注 -

Sun Cluster システムでは、次のクラスタファイルシステムをサポートします。

- Solaris ZFS™
- UNIX ファイルシステム (UFS)
- Sun StorEdge QFS ファイルシステム、および Sun QFS 共有ファイルシステム
- Sun Cluster プロキシファイルシステム (PxFS)
- Veritas ファイルシステム (VxFS)

HASStoragePlus リソースタイプは、確認、マウント、およびマウントの強制解除などの追加のファイルシステム機能を提供します。これらの機能により、Sun Cluster はローカルのファイルシステムをフェイルオーバーすることができます。フェイルオーバーを行うには、アフィニティスイッチオーバーが有効になった広域ディスクグループ上にローカルファイルシステムが存在していなければなりません。

HASStoragePlus リソースタイプの使用方法については、『[Sun Cluster データサービスの計画と管理 \(Solaris OS 版\)](#)』の「[高可用性ローカルファイルシステムの有効化](#)」を参照してください。

HASStoragePlus リソースタイプを使用して、リソースとリソースが依存するデバイスグループの起動を同期化することができます。詳細は、[84 ページ](#)の「[リソース、リソースグループ、リソースタイプ](#)」を参照してください。

## syncdir マウントオプション

syncdir マウントオプションは、実際に使用するファイルシステムとして UFS を使用するクラスタファイルシステムに使用できます。ただし、syncdir を指定しない場合、パフォーマンスは大幅に向上します。syncdir を指定した場合、POSIX 準拠の書き込みが保証されます。syncdir を指定しない場合、NFS ファイルシステムの場合と同じ動作となります。たとえば、syncdir を指定しないと、場合によっては、ファイルを閉じるまでスペース不足条件を検出できません。syncdir (および POSIX 動作) を指定すると、スペース不足条件は書き込み動作中に検出されます。syncdir を指定しない場合に問題が生じることはほとんどありません。

SPARC ベースのクラスタを使用している場合、VxFS には、UFS の syncdir マウントオプションと同等なマウントオプションはありません。VxFS の動作は syncdir マウントオプションを指定しない場合の UFS と同じです。

広域デバイスとクラスタファイルシステムの FAQ については、[112 ページ](#)の「[ファイルシステムに関する FAQ](#)」を参照してください。

# ディスクパスの監視

現在のリリースの Sun Cluster ソフトウェアは、ディスクパス監視機能 (DPM) をサポートします。この節では、DPM、DPM デーモン、およびディスクパスを監視するときに使用する管理ツールについての概念的な情報を説明します。ディスクパスの状態を監視、監視解除、および確認する手順については、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』を参照してください。

## DPM の概要

DPM は、二次ディスクパスの可用性を監視することによって、フェイルオーバーおよびスイッチオーバーの全体的な信頼性を向上させます。リソースを切り替える前には、`cldevice` コマンドを使用して、そのリソースが使用しているディスクパスの可用性を確認します。`cldevice` コマンドのオプションを使用すると、単一の Solaris ホストまたはクラスタ内のすべての Solaris ホストへのディスクパスを監視できます。コマンド行オプションの詳細は、[cldevice\(1CL\)](#) のマニュアルページを参照してください。

次の表に、DPM コンポーネントのデフォルトのインストール場所を示します。

| 場所                     | コンポーネント                                    |
|------------------------|--|
| デーモン                   | <code>/usr/cluster/lib/sc/scdpmd</code>    |
| コマンド行インタフェース           | <code>/usr/cluster/bin/cldevice</code>     |
| デーモン状態ファイル (実行時に作成される) | <code>/var/run/cluster/scdpm.status</code> |

マルチスレッド化された DPM デーモンは各ホスト上で動作します。DPM デーモン (`scdpmd`) はホストの起動時に `rc.d` スクリプトによって起動されます。問題が発生した場合、DPM デーモンは `pmfd` によって管理され、自動的に再起動されます。以下で、最初の起動時に `scdpmd` がどのように動作するかについて説明します。

注 - 起動時、各ディスクパスの状態は `UNKNOWN` に初期化されます。

1. DPM デーモンは、以前の状態ファイルまたは CCR データベースから、ディスクパスとノード名の情報を収集します。CCR についての詳細は、[50 ページの「クラスタ構成レポジトリ \(CCR\)」](#)を参照してください。DPM デーモンの起動後、指定したファイルから監視すべきディスクのリストを読み取るように DPM デーモンに指示できます。
2. DPM デーモンは通信インタフェースを初期化して、デーモンの外部にあるコンポーネント (コマンド行インタフェースなど) からの要求に応えます。

3. DPM デーモンは `scsi_inquiry` コマンドを使用して、監視リストにある各ディスクパスに 10 分ごとに `ping` を送信します。各エントリはロックされるため、通信インタフェースは監視中のエントリの内容にアクセスできなくなります。
4. DPM デーモンは、UNIX の `syslogd` コマンドを使用して Sun Cluster イベントフレームワークにパスの新しい状態を通知して、ログに記録します。詳細は、[syslogd\(1M\)](#) のマニュアルページを参照してください。

---

注 - このデーモンに関連するすべてのエラーは `pmfd` で報告されます。API のすべての関数は、成功時に `0` を返し、失敗時に `-1` を返します。

---

DPM デーモンは、Solaris I/O マルチパス (MPxIO) (従来の Sun StorEdge Traffic Manager)、Sun StorEdge 9900 Dynamic Link Manager、EMC PowerPath などのマルチパスドライバを通じて、論理パスの可用性を監視します。このようなマルチパスドライバは物理パスの障害を DPM デーモンから隠すため、DPM デーモンはマルチパスドライバが管理する物理パスを監視できません。

## ディスクパスの監視

この節では、クラスタ内のディスクパスを監視するための 2 つの方法について説明します。1 つめの方法は `cldevice` コマンドを使用する方法です。このコマンドを使用すると、クラスタ内のディスクパスの状態を監視、監視解除、または表示できます。このコマンドを使用して故障したディスクのリストを印刷し、ファイルからディスクパスを監視することもできます。詳細は、[cldevice\(1CL\)](#) のマニュアルページを参照してください。

2 つめの方法は、Sun Cluster Manager の GUI (Graphical User Interface) を使用してクラスタ内のディスクパスを監視する方法です。Sun Cluster Manager は、クラスタ内の監視しているディスクをトポロジビューで表示します。このトポロジビューは 10 分ごとに更新され、失敗した `ping` の数が表示されます。Sun Cluster Manager の GUI が報告する情報と `cldevice` コマンドを組み合わせると、ディスクパスを管理できます。Sun Cluster Manager については、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の第 12 章「グラフィカルユーザーインタフェースによる Sun Cluster の管理」を参照してください。

### cldevice コマンドを使用したディスクパスの監視と管理

`cldevice` コマンドを使用して、次の作業を実行できます。

- 新しいディスクパスの監視
- ディスクパスの監視解除
- CCR データベースからの構成データの再読み込み
- 指定したファイルからの監視または監視解除すべきディスクの読み取り
- クラスタ内の 1 つまたはすべてのディスクパスの状態の報告

- あるノードからアクセスできるすべてのディスクパスの印刷

任意のアクティブなノードから、ディスクパス引数を付けて `cldevice` コマンドを発行することによって、そのクラスタ上で DPM 管理作業を実行できます。ディスクパス引数はノード名とディスク名からなります。ノード名は不要です。ノード名を指定しない場合、デフォルトですべてのノードが影響を受けます。次の表に、ディスクパスの命名規約を示します。

注-必ず、UNIX のディスクパス名ではなく、広域ディスクパス名を指定してください。これは、広域ディスクのパス名がクラスタ全体で一貫しているためです。UNIX のディスクパス名にはこの性質はありません。たとえば、あるノードでディスクパス名を `c1t0d0` にして、別のノードで `c2t0d0` にすることができます。ノードに接続されたデバイスの広域ディスクパス名を調べるには、DPM コマンドを発行する前に `cldevice list` コマンドを使用します。詳細は、[cldevice\(1CL\)](#) のマニュアルページを参照してください。

表 3-3 ディスクパス名の例

| 名前型         | ディスクパス名の例                                 | 説明  |
|-------------|---|---|
| 広域ディスクパス    | <code>schost-1:/dev/did/dsk/d1</code>     | <code>schost-1</code> ノード上のディスクパス <code>d1</code>       |
|             | <code>all:d1</code>                       | クラスタのすべてのノードでのディスクパス <code>d1</code>                    |
| UNIX ディスクパス | <code>schost-1:/dev/rdisk/c0t0d0s0</code> | <code>schost-1</code> ノード上のディスクパス <code>c0t0d0s0</code> |
|             | <code>schost-1:all</code>                 | <code>schost-1</code> ノードでのすべてのディスクパス                   |
| すべてのディスクパス  | <code>all:all</code>                      | クラスタのすべてのノードでのすべてのディスクパス                                |

## Sun Cluster Manager によるディスクパスの監視

Sun Cluster Manager を使用すると、次のような DPM の基本的な管理作業を実行できます。

- ディスクパスの監視
- ディスクパスの監視解除
- クラスタ内のすべての監視対象ディスクパスの状態の表示
- 監視されているすべてのディスクパスが失敗したときの Solaris ホストの自動再起動の有効化または無効化

Sun Cluster Manager のオンラインヘルプでは、ディスクパスの管理方法の手順について説明しています。

`clnode set` コマンドを使用して、ディスクパスエラーを管理するすべての監視対象ディスクパスでエラーが発生した際のノードの自動再起動を有効化または無効化するには、`clnode set` コマンドを使用します。Sun Cluster Manager を使用してこれらの作業を実行することもできます。

## 定足数と定足数デバイス

この節には、次のトピックが含まれます。

- 64 ページの「定足数投票数について」
- 65 ページの「定足数の構成について」
- 65 ページの「定足数デバイス要件の順守」
- 66 ページの「定足数デバイスのベストプラクティスの順守」
- 67 ページの「推奨される定足数の構成」
- 69 ページの「変則的な定足数の構成」
- 70 ページの「望ましくない定足数の構成」

---

注 - Sun Cluster ソフトウェアが定足数デバイスとしてサポートする特定のデバイスの一覧については、Sun のサービスプロバイダにお問い合わせください。

---

クラスタノードはデータとリソースを共有しており、複数のアクティブなパーティションがあるとデータが壊れる恐れがあるのでクラスタは決して複数のアクティブなパーティションに一度に分割しないでください。クラスタメンバーシップモニター (Cluster Membership Monitor、CMM) および定足数アルゴリズムにより、たとえクラスタ接続がパーティション分割されている場合でも、いつでも同じクラスタのインスタンスが 1 つだけは動作していることが保証されます。

定足数と CMM の概要については、『[Sun Cluster の概要 \(Solaris OS 版\)](#)』の「[クラスタメンバーシップ](#)」を参照してください。

クラスタのパーティション分割からは、次の 2 種類の問題が発生します。

- Split brain
- amnesia

*Split brain* は、ノード間のクラスタ接続が失われ、クラスタがサブクラスタにパーティション分割されるときに起きます。あるパーティションのノードはほかのパーティションのノードと通信できないため、各パーティションは自分が唯一のパーティションであると認識します。

*amnesia* は、停止したクラスタが、停止時よりも古いクラスタ構成データに基づいて再起動されたときに発生します。この問題は、最後に機能していたクラスタパーティションにないノード上のクラスタを起動するときに起きる可能性があります。

Sun Cluster ソフトウェアは、*split brain* と *amnesia* を次の操作により回避します。

- 各ノードに1つの投票を割り当てる
- 動作中のクラスタの過半数の投票を管理する

過半数の投票数を持つパーティションは、定足数を獲得し、動作可能になります。この過半数の投票メカニズムにより、クラスタ内に3つ以上のノードが構成されているときに *split brain* と *amnesia* を防ぐことができます。ただし、クラスタ内に3つ以上のノードが構成されている場合、ノードの投票数を数えるだけでは十分ではありません。しかし、2ホストクラスタでは過半数が2であるため、このような2ホストクラスタがパーティション分割された場合、いずれかのパーティションが定足数を獲得するために外部投票が必要です。この外部からの投票は「定足数デバイス」によって行われます。

## 定足数投票数について

`clquorum show` コマンドを使って、次の情報を調べます。

- 構成済み投票数
- 現在の投票数
- 定足数に必要な投票数

詳細は、`cluster(1CL)` のマニュアルページを参照してください。

ノードおよび定足数デバイスの両方がクラスタへの投票に数えられ、定足数を満たすことができます。

ノードは、ノードの状態に応じて投票に数えられます。

- ノードが起動してクラスタメンバーになると、投票数は1となります。
- ノードがインストールされているときは、投票数は0となります。
- システム管理者がノードを保守状態にすると、投票数は0となります。

定足数デバイスは、デバイスに伴う投票数に基づいて、投票に数えられます。定足数デバイスを構成するとき、Sun Cluster ソフトウェアは定足数デバイスに  $N-1$  の投票数を割り当てます ( $N$  は定足数デバイスに伴う投票数)。たとえば、2つのノードに接続された、投票数がゼロ以外の定足数デバイスの投票数は  $1(2-1)$  になります。



定足数デバイスは、次の2つの条件のうちの1つを満たす場合に投票に数えられます。

- 定足数デバイスに現在接続されている1つ以上のノードがクラスタメンバーである。
- 定足数デバイスに現在接続されている1つ以上のホストが起動中で、そのホストは定足数デバイスを所有する最後のクラスタパーティションのメンバーであった。

定足数デバイスの構成は、クラスタをインストールするときに行うか、あるいはあとで、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の第6章「定足数の管理」で説明されている手順に従って行います。

## 定足数の構成について

次に、定足数の構成について示します。

- 定足数デバイスには、ユーザーデータを含むことができます。
- $N+1$  の構成 ( $N$ 個の定足数デバイスがそれぞれ、1から  $N$ までの Solaris ホストのうちの1つのホストと  $N+1$  番目の Solaris ホストに接続されている構成) では、1から  $N$ までのどの Solaris ホストで障害が発生しても、 $N/2$  個のうちの任意の Solaris ホストに障害が発生しても、そのクラスタは影響を受けません。この可用性は、定足数デバイスが正しく機能していることを前提にしています。
- $N$ ホスト構成 (1つの定足数デバイスがすべてのホストに接続されている構成) では、 $N-1$  個のうちの任意のホストに障害が発生しても、そのクラスタは影響を受けません。この可用性は、定足数デバイスが正しく機能していることを前提にしています。
- 1つの定足数デバイスがすべてのホストに接続している  $N$ ホスト構成では、すべてのクラスタホストが使用できる場合、定足数デバイスに障害が起きてもクラスタは影響を受けません。

回避すべき定足数の構成例については、[70 ページ](#)の「望ましくない定足数の構成」を参照してください。推奨される定足数の構成例については、[67 ページ](#)の「推奨される定足数の構成」を参照してください。

## 定足数デバイス要件の順守

Sun Cluster ソフトウェアがご使用のデバイスを定足数デバイスとしてサポートしていることを確認します。この要件を無視すると、クラスタの可用性が損なわれる場合があります。

---

注 - Sun Cluster ソフトウェアが定足数デバイスとしてサポートする特定のデバイスの一覧については、Sun のサービスプロバイダにお問い合わせください。

---

Sun Cluster ソフトウェアは、次の種類の定足数デバイスをサポートしています。

- SCSI-3 PGR 予約に対応した多重ホスト共有ディスク。
- SCSI-2 予約に対応した二重ホスト共有ディスク。
- Sun または Network Appliance, Incorporated の NAS (Network-Attached Storage) デバイス。
- 定足数サーバーマシン上で動作する定足数サーバープロセス。
- ディスクのフェンシングをオフに切り替え、ソフトウェア定足数を使用している場合は任意の共有ディスク。ソフトウェア定足数は、Sun が開発したプロトコルで、SCSI Persistent Group Reservations (PGR) の形成をエミュレートします。



---

注意 - SATA (Serial Advanced Technology Attachment) ディスクなど、SCSI に対応していないディスクを使用する場合は、フェンシングをオフにしてください。

---

---

注 - レプリケートされたデバイスは定足数デバイスとして使用できません。

---

2 ホスト構成では、1 つのホストに障害が起きてももう 1 つのホストが動作を継続できるように、少なくとも 1 つの定足数デバイスを構成する必要があります。詳細は、[図 3-2](#) を参照してください。

回避すべき定足数の構成例については、[70 ページ](#)の「望ましくない定足数の構成」を参照してください。推奨される定足数の構成例については、[67 ページ](#)の「推奨される定足数の構成」を参照してください。

## 定足数デバイスのベストプラクティスの順守

以下の情報を使用して、ご使用のトポロジに最適な定足数の構成を評価してください。

- クラスタの全 Solaris ホストに接続できるデバイスがありますか。
  - ある場合は、そのデバイスを 1 つの定足数デバイスとして構成してください。この構成は最適な構成なので、別の定足数デバイスを構成する必要はありません。



---

注意-この要件を無視して別の定足数デバイスを追加すると、追加した定足数デバイスによってクラスタの可用性が低下します。

---

- ない場合は、1つまたは複数のデュアルポートデバイスを構成してください。
- 定足数デバイスにより提供される投票の合計数が、ノードにより提供される投票の合計数より必ず少なくなるようにします。少なくなければ、すべてのノードが機能していても、すべてのディスクを使用できない場合、そのノードはクラスタを形成できません。

---

注-特定の環境によっては、自分のニーズに合うように、全体的なクラスタの可用性を低くした方が望ましい場合があります。このような場合には、このベストプラクティスを無視できます。ただし、このベストプラクティスを守らないと、全体の可用性が低下します。たとえば、[69 ページの「変則的な定足数の構成」](#)に記載されている構成では、クラスタの可用性は低下し、定足数の投票がノードの投票を上回ります。クラスタでは、Host A と Host B 間にある共有ストレージへのアクセスが失われると、クラスタ全体に障害が発生します。

---

このベストプラクティスの例外については、[69 ページの「変則的な定足数の構成」](#)を参照してください。

- 記憶装置へのアクセスを共有するホストのすべてのペア間で定足数デバイスを指定します。この定足数の構成により、障害からの影響の防止プロセスが高速化されます。詳細は、[68 ページの「2 ホストより大きな構成での定足数」](#)を参照してください。
- 通常、定足数デバイスの追加によりクラスタの投票の合計数が同じになる場合、クラスタ全体の可用性は低下します。
- ノードを追加したり、ノードに障害が発生すると、定足数デバイスの再構成は少し遅くなります。従って、必要以上の定足数デバイスを追加しないでください。

回避すべき定足数の構成例については、[70 ページの「望ましくない定足数の構成」](#)を参照してください。推奨される定足数の構成例については、[67 ページの「推奨される定足数の構成」](#)を参照してください。

## 推奨される定足数の構成

この節では、推奨される定足数の構成例を示します。回避すべき定足数の構成例については、[70 ページの「望ましくない定足数の構成」](#)を参照してください。

## 2 ホスト構成の定足数

2ホストのクラスタを形成するには、2つの定足投票数が必要です。これらの2つの投票数は、2つのクラスタホスト、または1つのホストと1つの定足数デバイスのどちらかによるものです。

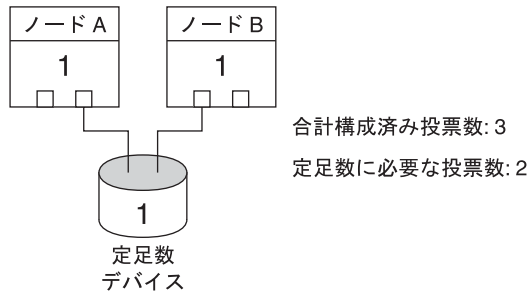
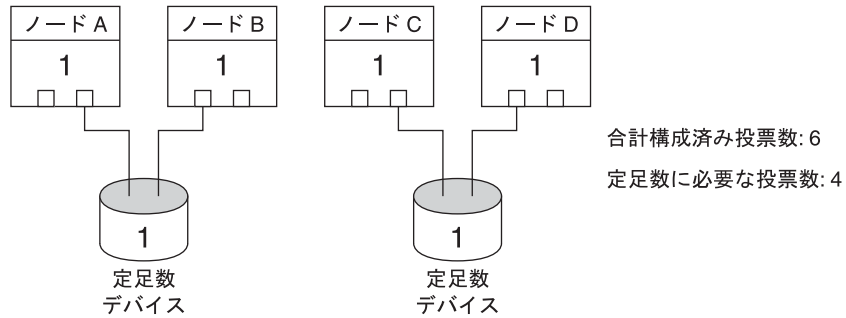


図 3-2 2ホスト構成

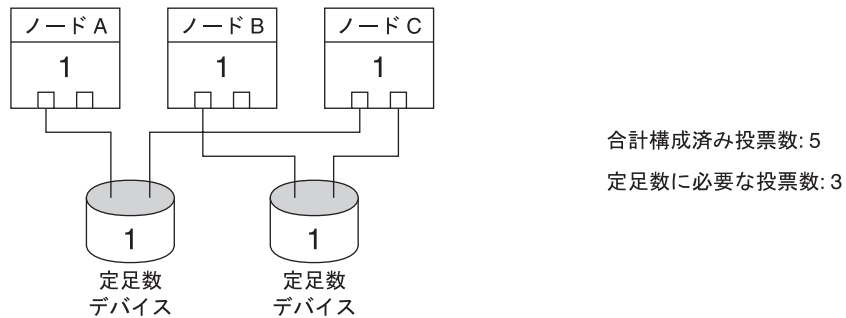
## 2ホストより大きな構成での定足数

クラスタに3つ以上のホストが含まれている場合、定足数デバイスを持たない1つのホストに障害が発生してもクラスタはその影響を受けないので、定足数デバイスは必要ありません。ただし、このような条件下では、クラスタ内に過半数のホストがないとクラスタを起動できません。

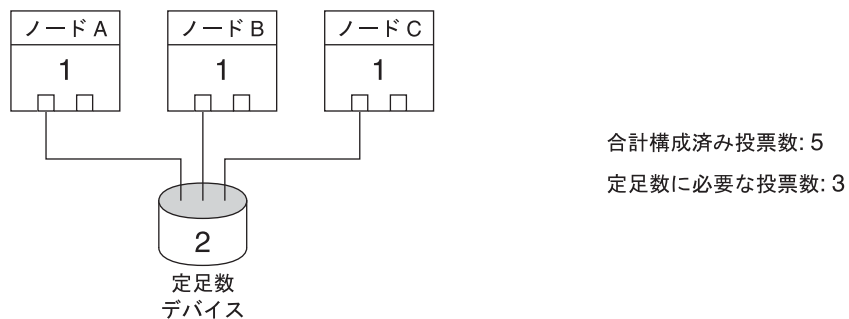
3つ以上のホストを含むクラスタに定足数デバイスを追加できます。ホストおよび定足数デバイスの投票数を含め、パーティションが定足投票数の過半数を獲得している場合、そのパーティションはクラスタとして存続できます。したがって、定足数デバイスを追加する場合は、定足数デバイスを構成するかどうか、および定足数デバイスの構成先を選択するときに、発生する可能性があるホストおよび定足数デバイスの障害を考慮するようにしてください。



この構成は、どちらかのペアが稼働し続けるためには各ペアが稼働していなければならない。



この構成は、通常、アプリケーションがノード A とノード B で実行されるように構成され、ノード C をホットスペアとして使用する。



この構成は、任意の1つ以上のノードと定足数デバイスとの組み合わせでクラスタを形成できる。

## 変則的な定足数の構成

図 3-3 では、Host A と Host B でミッションクリティカルなアプリケーション (Oracle データベースなど) を実行していると仮定します。Host A と Host B を使用できず、共

有データにアクセスできない場合、クラスタ全体を停止させる必要がある場合があります。停止させない場合、この構成は高可用性を提供できないため、最適な構成とはなりません。

この例外に関するベストプラクティスについては、66ページの「定足数デバイスのベストプラクティスの順守」を参照してください。

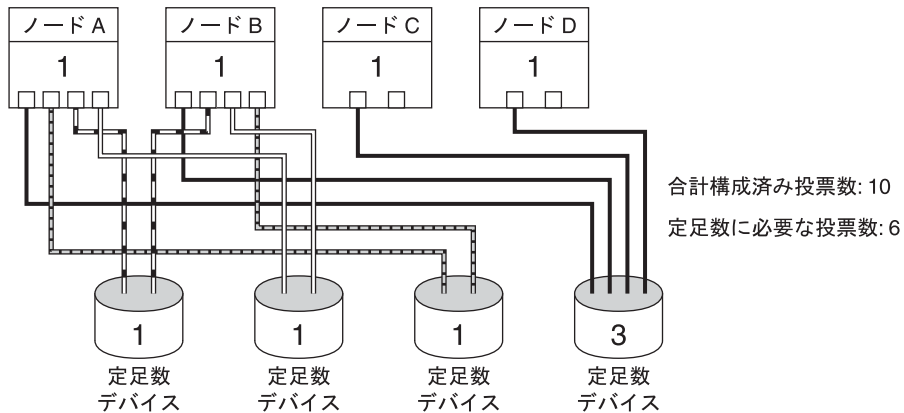
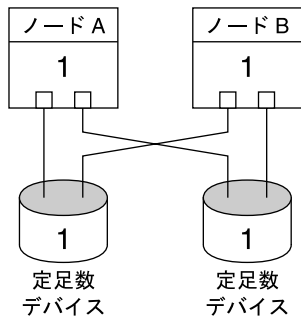


図 3-3 変則的な構成

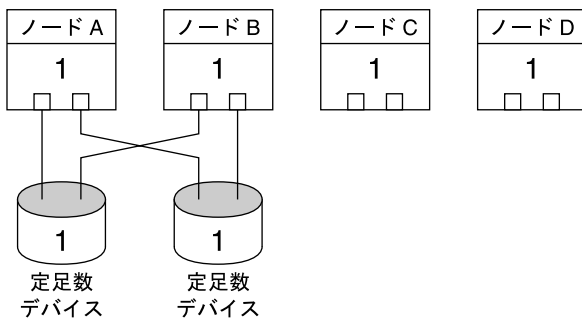
## 望ましくない定足数の構成

この節では、回避すべき定足数の構成例を示します。推奨される定足数の構成例については、67ページの「推奨される定足数の構成」を参照してください。



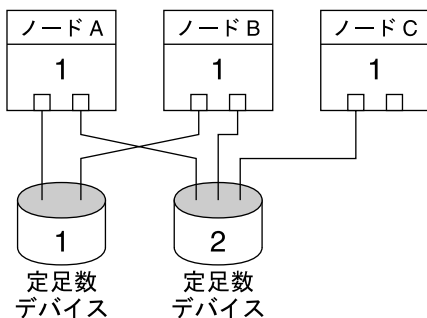
合計構成済み投票数: 4  
定足数に必要な投票数: 3

この構成は、定足数デバイス投票が、ノードの投票数より少なくなければならぬというベストプラクティスに反する。



合計構成済み投票数: 6  
定足数に必要な投票数: 4

この構成は、定足数デバイスを追加して、合計投票数を等しくするべきではないというベストプラクティスに反する。この構成では、可用性は向上されない。



合計構成済み投票数: 6  
定足数に必要な投票数: 4

この構成は、定足数デバイス投票が、ノードの投票数より少なくなければならぬというベストプラクティスに反する。

## データサービス

「データサービス」という用語は、Oracle や Sun Java System Web Server など、単一のサーバーではなく、クラスタで動作するように構成されたアプリケーションを意味

します。データサービスは、アプリケーション、専用の Sun Cluster 構成ファイル、および、アプリケーションの次のアクションを制御する Sun Cluster 管理メソッドからなります。

- 開始
- 停止
- 監視と訂正手段の実行

データサービスタイプについては、『Sun Cluster の概要 (Solaris OS 版)』の「データサービス」を参照してください。

図 3-4 に、単一のアプリケーションサーバーで動作するアプリケーション (単一サーバーモデル) と、クラスタで動作する同じアプリケーション (クラスタサーバーモデル) との比較を示します。これら 2 つの構成の唯一の違いは、クラスタ化されたアプリケーションの動作がより速くなり、その可用性もより高くなることです。

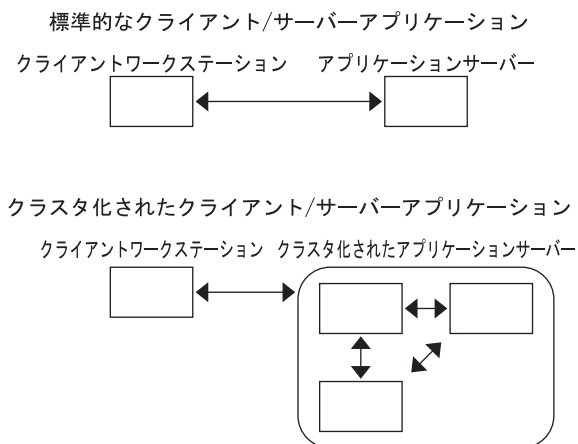


図 3-4 標準的なクライアントサーバー構成とクラスタ化されたクライアントサーバー構成

単一サーバーモデルでは、特定のパブリックネットワークインタフェース (ホスト名) を介してサーバーにアクセスするようにアプリケーションを構成します。ホスト名はその物理サーバーに関連付けられています。

クラスタサーバーモデルでは、パブリックネットワークインタフェースは、論理ホスト名または共有アドレスです。「ネットワークリソース」は、論理ホスト名と共有アドレスの両方を表します。

一部のデータサービスでは、ネットワークインタフェースとして論理ホスト名か共有アドレスのいずれかを指定する必要があります。論理ホスト名と共有アドレスは相互に交換できません。しかし、別のデータサービスでは、論理ホスト名や共有ア



ドレスをどちらでも指定することができます。どのようなタイプのインタフェースを指定する必要があるかどうかについては、各データサービスのインストールや構成の資料を参照してください。

ネットワークリソースは、特定の物理サーバーと関連付けられているわけではありません。ネットワークリソースは、ある物理サーバーから別の物理サーバーに移すことができます。

ネットワークリソースは、当初、1つのノード(主ノード)に関連付けられています。主ノードで障害が発生すると、ネットワークリソースとアプリケーションリソースは別のクラスタノード(二次ノード)にフェイルオーバーされます。ネットワークリソースがフェイルオーバーされても、アプリケーションリソースは、短時間の遅れの後に二次ノードで動作を続けます。

図3-5に、単一サーバーモデルとクラスタサーバーモデルとの比較を示します。クラスタサーバーモデルのネットワークリソース(この例では論理ホスト名)は、複数のクラスタノード間を移動できます。アプリケーションは、特定のサーバーに関連付けられたホスト名として、この論理ホスト名を使用するように設定されます。

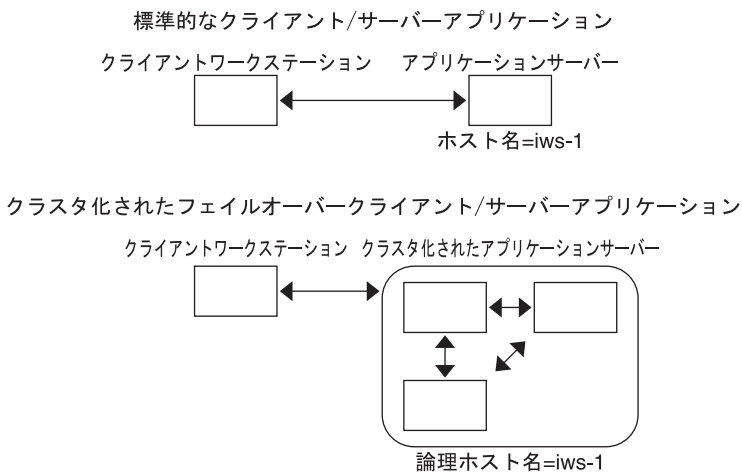


図3-5 固定ホスト名と論理ホスト名

共有アドレスも最初は1つのノードに対応付けられます。このノードのことを「広域インタフェースノード」といいます。共有アドレスは「広域インタフェース」といい、クラスタへの単一ネットワークインタフェースとして使用されます。

論理ホスト名モデルとスケーラブルサービスモデルの違いは、スケーラブルサービスモデルでは、各ノードのループバックインタフェースにも共有アドレスがアクティブに構成される点です。この構成では、データサービスの複数のインスタンスをいくつかのノードで同時にアクティブにすることができます。「スケーラブルサー

ビス」という用語は、クラスタノードを追加してアプリケーションのCPUパワーを強化すれば、性能が向上することを意味します。

広域インタフェースノードに障害が発生した場合、共有アドレスは同じアプリケーションのインスタンスが動作している別のノードで起動できます。これによって、このノードが新しい広域インタフェースノードになります。または、共有アドレスを、このアプリケーションを実行していない別のクラスタノードにフェイルオーバーさせることができます。

図3-6に、単一サーバー構成とクラスタ化されたスケーラブルサービス構成との比較を示します。スケーラブルサービス構成では、共有アドレスがすべてのノードに設定されています。アプリケーションは、特定のサーバーに関連付けられたホスト名として、この共有アドレスを使用するように設定されます。この方式は、フェイルオーバーデータサービスに論理ホスト名を使用する方法と似ています。

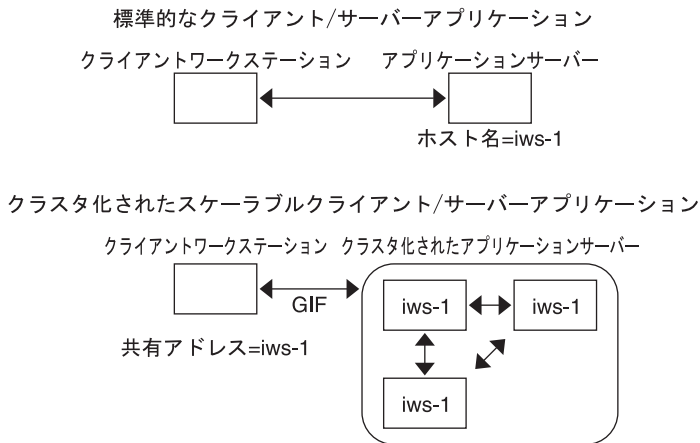


図3-6 固定ホスト名と共有アドレス

## データサービスメソッド

Sun Cluster ソフトウェアでは、一連のサービス管理メソッドを提供しています。これらのメソッドは、リソースグループマネージャー (Resource Group Manager、RGM) の制御下で実行されます。RGMはこれらのメソッドを使用して、クラスタノード上のアプリケーションを開始、停止、および監視します。これらのメソッドとクラスタフレームワークソフトウェアおよび多重ホストデバイスにより、アプリケーションは、フェイルオーバーデータサービスやスケーラブルデータサービスとして機能します。

さらに、RGMは、アプリケーションのインスタンスやネットワークリソース (論理ホスト名と共有アドレス) といったクラスタのリソースを管理します。

Sun Cluster ソフトウェアが提供するメソッドのほかにも、Sun Cluster ソフトウェアは API やいくつかのデータサービス開発ツールを提供します。これらのツールを使用すれば、アプリケーション開発者は、独自のデータサービスメソッドを開発することによって、ほかのアプリケーションを Sun Cluster ソフトウェアの下で高可用性データサービスとして実行できます。

## フェイルオーバーデータサービス

データサービスが実行されているノード(主ノード)に障害が発生すると、サービスは、ユーザーによる介入なしで別の作業ノードに移行します。フェイルオーバーサービスは、アプリケーションインスタンスリソースとネットワークリソース(「論理ホスト名」)のコンテナである「フェイルオーバーリソースグループ」を使用します。論理ホスト名は、1つのノードで構成でき、あとで自動的に元のノードに構成したり、別のノードに構成したりすることができる IP アドレスです。

フェイルオーバーデータサービスでは、アプリケーションインスタンスは単一のノードでのみ実行されます。フォルトモニターは、エラーを検出すると、そのインスタンスを同じノードで再起動しようとするか、別のノードで起動(フェイルオーバー)しようとしています。出力は、データサービスの構成方法によって異なります。

## スケーラブルデータサービス

スケーラブルデータサービスは、複数のノードのアクティブインスタンスに対して効果があります。

スケーラブルサービスは、2つのリソースグループを使用します。

- 「スケーラブルリソースグループ」には、アプリケーションリソースが含まれません。
- フェイルオーバーリソースグループには、スケーラブルサービスが依存するネットワークリソース(「共有アドレス」)が含まれます。共有アドレスは、ネットワークアドレスです。このネットワークアドレスは、クラスタ内のノードのすべてのスケーラブルサービスによってバインドできます。この共有アドレスにより、これらのスケーラブルサービスをノードに拡張できます。クラスタには複数の共有アドレスを設定することができ、1つのサービスは、複数の共有アドレスにバインドすることができます。

スケーラブルリソースグループは、同時に複数のノードでオンラインにできます。その結果、サービスの複数のインスタンスを一度に実行できます。スケーラブルなリソースグループはすべて負荷分散を使用します。スケーラブルサービスをホストとするすべてのノードは、サービスをホストするための同じ共有アドレスを使用します。共有アドレスのホストとなるフェイルオーバーリソースグループは、一度に1つのノードでしかオンラインにできません。

サービス要求は、単一ネットワークインタフェース (広域インタフェース) を通じてクラスタに入ります。これらの要求は、事前に定義されたいくつかのアルゴリズムの1つに基づいてノードに分配されます。これらのアルゴリズムは「負荷均衡ポリシー」によって設定されます。クラスタは、負荷均衡ポリシーを使用し、いくつかのノード間でサービス負荷均衡をとることができます。ほかの共有アドレスをホストしている異なるノード上には、複数の広域インタフェースが存在する可能性があります。

スケラブルサービスの場合、アプリケーションインスタンスはいくつかのノードで同時に実行されます。広域インタフェースのホストとなるノードに障害が発生すると、広域インタフェースは別のノードで処理を続行します。動作していたアプリケーションインスタンスが停止した場合、そのインスタンスは同じノード上で再起動しようとしています。

アプリケーションインスタンスを同じノードで再起動できず、別の未使用のノードがサービスを実行するように構成されている場合、サービスはその未使用ノードで処理を続行します。そうしないと、サービスは残りのノード上で動作し続け、サービスのスループットが低下することがあります。

---

注 - 各アプリケーションインスタンスの TCP 状態は、広域インタフェースノードではなく、インスタンスを持つノードで維持されます。したがって、広域インタフェースノードに障害が発生しても接続には影響しません。

---

図 3-7 に、フェイルオーバーリソースグループとスケラブルリソースグループの例と、スケラブルサービスにとってそれらの間にどのような依存関係があるのかを示します。この例は、3つのリソースグループを示しています。フェイルオーバーリソースグループには、可用性の高い DNS のアプリケーションリソースと、可用性の高い DNS および可用性の高い Apache Web Server (SPARC ベースのクラスタに限り使用可能) の両方から使用されるネットワークリソースが含まれます。スケラブルリソースグループには、Apache Web Server のアプリケーションインスタンスだけが含まれます。リソースグループの依存関係は、スケラブルリソースグループとフェイルオーバーリソースグループの間に存在します (実線)。また、Apache アプリケーションリソースはすべて、共有アドレスであるネットワークリソース schost-2 に依存します (破線)。

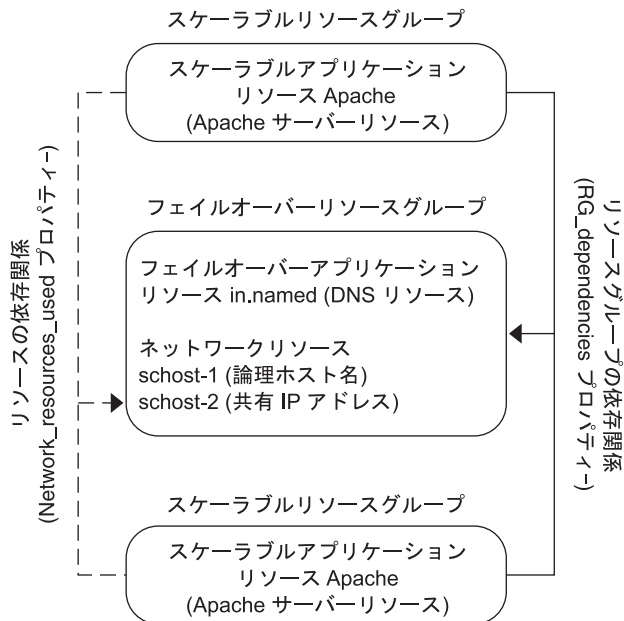


図 3-7 SPARC: フェイルオーバーリソースグループとスケーラブルリソースグループの例

## 負荷均衡ポリシー

負荷均衡は、スケーラブルサービスのパフォーマンスを応答時間とスループットの両方の点で向上させます。スケーラブルデータサービスには2つのクラスがあります。

- pure
- sticky

*pure* サービスでは、任意のインスタンスがクライアント要求に応答できます。*sticky* サービスでは、クライアントは同じインスタンスに応答を送信できます。これらの要求は、別のインスタンスには変更されません。

*pure* サービスは、ウェイト設定した (*weighted*) 負荷均衡ポリシーを使用します。この負荷均衡ポリシーのもとでは、クライアント要求は、デフォルトで、クラスタ内のサーバーインスタンスに一律に分配されます。負荷は指定されたウェイト値に従って各種のノードに分配されます。たとえば、3 ノードクラスタにおいて、各ノードのウェイトが1だとします。各ノードは、そのサービスに対する任意のクライアントからの要求を 1/3 ずつ負担します。クラスタ管理者は、管理コマンドまたは Sun Cluster Manager によって、いつでもウェイト値を変更できます。

ウェイト設定した負荷分散ポリシーは、`Load_balancing_weights` プロパティに設定された `LB_WEIGHTED` 値を使用して設定されます。ウェイトがノードについて明示的に設定されていない場合は、デフォルトで1が設定されます。

ウェイト設定したポリシーは、一定の割合のクライアントトラフィックを特定ノードに送るためのものです。たとえば、 $X$  = 「ウェイト」、 $A$  = 「すべてのアクティブノードの合計ウェイト」であるとします。アクティブノードは、新しい接続数の合計の約  $X/A$  がこのアクティブノード宛てに送られると予測できます。ただし、接続数の合計は十分に大きな数である必要があります。このポリシーは、個々の要求には対応しません。

このウェイト設定したポリシーは、ラウンドロビンではないことに注意してください。ラウンドロビンポリシーでは、常に、同じクライアントからの要求はそれぞれ異なるノードに送られます。たとえば、1 番目の要求はノード 1 に、2 番目の要求はノード 2 に送られます。

sticky サービスには「*ordinary sticky*」と「*wildcard sticky*」の 2 種類があります。

sticky サービスを使用すると、内部状態メモリーを共有でき (アプリケーションセッション状態)、複数の TCP 接続でアプリケーションレベルの同時セッションが可能です。

*ordinary sticky* サービスを使用すると、クライアントは、複数の同時 TCP 接続で状態を共有できます。このクライアントを、単一ポートで待機するサーバーインスタンスに対して「sticky」であるといいます。

次の条件を満たす場合、このクライアントはすべての要求が同じサーバーインスタンスに送られることが保証されます。

- インスタンスが動作中でアクセス可能であること。
- サービスがオンラインの間、負荷分散ポリシーが変更されないこと。

たとえば、クライアント上の Web ブラウザは、3 つの異なる TCP 接続を使用して、ポート 80 にある共有 IP アドレスに接続します。そして、これらの接続はサービスでキャッシュされたセッション情報をお互いに交換します。

sticky ポリシーを一般化すると、そのポリシーは同じインスタンスの背後でセッション情報を交換する複数のスケーラブルサービスにまで及びます。これらのサービスが同じインスタンスの背後でセッション情報を交換するとき、同じノードで異なるポートと通信する複数のサーバーインスタンスに対して、そのクライアントは「sticky」であるといいます。

たとえば、電子商取引 Web サイトの顧客はポート 80 の HTTP を使用して買い物をします。そして、購入した製品をクレジットカードで支払うときには、ポート 443 の SSL に切り替えて機密データを送ります。

通常の sticky ポリシーでは、ポートの集合が、アプリケーションリソースの構成時に認識されます。このポリシーは、`Load_balancing_policy` リソースプロパティの `LB_STICKY` の値を使用して設定されます。

*Wildcard sticky* サービスは、動的に割り当てられたポート番号を使用しますが、クライアント要求が同じノードに送りかえされると想定します。クライアントは、同じ IP アドレスを持っているポートに対して「sticky wildcard」であるといいます。

このポリシーの例としては、受動モード FTP があります。たとえば、クライアントはポート 21 の FTP サーバーに接続します。次に、サーバーは、動的ポート範囲のリスナーポートサーバーに接続し直すようにクライアントに指示します。この IP アドレスに対する要求はすべて、サーバーが制御情報によってクライアントに通知した、同じノードに転送されます。

sticky-wildcard ポリシーは、通常の「sticky」ポリシーの上位セットです。IP アドレスによって識別されるスケーラブルサービスでは、ポートはサーバーによって割り当てられます。したがって、事前に認識できません。ポートは変更されることがあります。このポリシーは、Load\_balancing\_policy リソースプロパティの LB\_STICKY\_WILD の値を使用して設定されます。

このような各 sticky ポリシーでは、ウェイト設定した負荷均衡ポリシーがデフォルトで有効です。したがって、クライアントの最初の要求は、負荷均衡によって指定されたインスタンス宛てに送られます。インスタンスが動作しているノードとのアフィニティーをクライアントが確立すると、今後の要求はそのインスタンス宛てに送られます。ただし、そのノードはアクセス可能であり、負荷均衡ポリシーが変更されていない必要があります。

## フェイルバック設定

リソースグループは、あるノードから別のノードへ処理を継続します。このようなフェイルオーバーが発生すると、二次ノードが新しい主ノードになります。フェイルバック設定は、本来の主ノードがオンラインに戻ったときの動作を指定します。つまり、本来の主ノードを再び主ノードにする(フェイルバックする)か、現在的主ノードをそのままにするかです。これを指定するには、Failback リソースグループプロパティ設定を使用します。

リソースグループをホストしていた本来の主ノードに障害が発生し、繰り返し再起動する場合は、フェイルバックを設定することによって、リソースグループの可用性が低くなる可能性もあります。

## データサービス障害モニター

Sun Cluster の各データサービスには、データサービスを定期的に検査してその状態を判断するフォルトモニターがあります。フォルトモニターは、アプリケーションデーモンが動作しているかどうかや、クライアントにサービスが提供されているかどうかを検証します。探索によって得られた情報をもとに、デーモンの再起動やフェイルオーバーの実行などの事前に定義された処置が開始されます。

## 新しいデータサービスの開発

Sun が提供する構成ファイルや管理メソッドのテンプレートを使用することで、さまざまなアプリケーションをクラスタ内でフェイルオーバーサービスやスケラブルサービスとして実行できます。フェイルオーバーサービスやスケラブルサービスとして実行するアプリケーションが Sun から提供されていない場合は、代替の方法があります。つまり、Sun Cluster API や DSET API を使用して、フェイルオーバーサービスやスケラブルサービスとして動作するようにアプリケーションを構成します。しかし、必ずしもすべてのアプリケーションがスケラブルサービスになるわけではありません。

### スケラブルサービスの特徴

アプリケーションがスケラブルサービスになれるかどうかを判断するには、いくつかの基準があります。アプリケーションがスケラブルサービスになれるかどうかを判断する方法については、『[Sun Cluster データサービス開発ガイド \(Solaris OS 版\)](#)』の「[アプリケーションの適合性の分析](#)」を参照してください。



次に、これらの基準の要約を示します。

- 第1に、このようなサービスは1つまたは複数のサーバーインスタンスから構成されます。各インスタンスは、異なるノードで実行されます。同じサービスの複数のインスタンスを、同じノードで実行することはできません。
- 第2に、このサービスが外部の論理データストアを提供する場合は、十分に注意する必要があります。このストアに複数のサーバーインスタンスから並行にアクセスする場合、このような並行アクセスの同期をとることによって、更新を失ったり、変更中のデータを読み取ったりすることを避ける必要があります。「外部」という用語は、このストアとメモリー内の状態を区別するために使用しています。「論理」という用語は、ストア自体複製されている場合でも、単一の実体として見えることを表します。さらに、このデータストアでは、サーバーインスタンスがデータストアを更新すると、その更新がすぐにほかのインスタンスで「認識」されます。

Sun Cluster ソフトウェアは、このような外部記憶領域をそのクラスタファイルシステムと広域 raw パーティションを介して提供します。例として、サービスが外部ログファイルに新しいデータを書き込む場合や既存のデータを修正する場合を想定してください。このサービスのインスタンスが複数実行されている場合、各インスタンスはこの外部ログへのアクセスを持ち、このログに同時にアクセスできます。各インスタンスは、このログに対するアクセスの同期をとる必要があります。そうしないと、インスタンスは相互に妨害しあうことになります。このサービスは、`fcntl` と `lockf` による通常の Solaris ファイルロック機能を使用して、必要な同期をとることができます。

この種類のデータストアのもう一つの例はバックエンドデータベースで、たとえば、Oracle や SPARC ベースのクラスタ用の高可用性 Oracle Real Application Clusters Guard などがあります。この種類のバックエンドデータベースサーバーには、データベース照会または更新トランザクションによる同期が組み込まれています。したがって、複数のサーバーインスタンスが独自の同期を実装する必要はありません。

スケラブルサービスではないサービスの例としては、Sun の IMAP サーバーがあります。このサービスは記憶領域を更新しますが、その記憶領域はプライベートであり、複数の IMAP インスタンスがこの記憶領域に書き込むと、更新の同期がとられないために相互に上書きし合うことになります。IMAP サーバーは、同時アクセスの同期をとるよう書き直す必要があります。

- 最後に、インスタンスは、ほかのインスタンスのデータから切り離されたプライベートデータを持つ場合があることに注意してください。このような場合、データはプライベートであり、そのインスタンスだけがデータを処理するため、サービスは並行アクセスの同期をとる必要はありません。この場合、個人的なデータをクラスタファイルシステムに保存すると、これらのデータが広域にアクセス可能になる可能性があるため、十分に注意する必要があります。

## データサービス API と DSDL API

Sun Cluster ソフトウェアには、アプリケーションの可用性を高めるものとして次の機能があります。

- Sun Cluster ソフトウェアの一部として提供されるデータサービス
- データサービス API
- データサービス用の開発ライブラリ API
- 汎用データサービス

Sun Cluster ソフトウェアが提供するデータサービスをインストールおよび構成する方法については、『[Sun Cluster データサービスの計画と管理 \(Solaris OS 版\)](#)』を参照してください。Sun Cluster フレームワークでアプリケーションの可用性を高める方法については、Sun Cluster 3.1 9/04 Software Collection for Solaris OS (SPARC Platform Edition)を参照してください。

アプリケーション開発者は、Sun Cluster API を使用することによって、データサービスインスタンスの起動や停止を行なう障害モニターやスクリプトを開発できます。これらのツールを使用すると、アプリケーションをフェイルオーバーまたはスケラブルデータサービスとして実装できます。Sun Cluster ソフトウェアは「汎用」のデータサービスを提供します。この汎用のデータサービスを使用すれば、簡単に、アプリケーションに必要な起動メソッドと停止メソッドを生成して、データサービスをフェイルオーバーサービスまたはスケラブルサービスとして実装できます。

## クラスタインターコネクトによるデータサービストラフィックの送受信

通常、クラスタには Solaris ホスト間を結ぶ複数のネットワーク接続が必要です。クラスタインターコネクトは、これらの接続から構成されています。

Sun Cluster ソフトウェアは、複数のインターコネクトを使用して、次の目標を達成します。

- 高可用性を保証します。
- 性能を向上させます。

内部と外部の両方のトラフィック (ファイルシステムのデータやスケラブルサービスのデータなど) の場合、メッセージはすべての利用できるインターコネクト間で転送されます。クラスタインターコネクトは、ホスト間の通信の可用性を高めるためにアプリケーションから使用することもできます。たとえば、分散アプリケーションでは、個々のコンポーネントが異なるホストで動作することがあり、その場合には、ホスト間の通信が必要になります。パブリック伝送の代わりにクラスタインターコネクトを使用することで、個別のリンクに障害が発生しても、接続を持続することができます。

ホスト間の通信にクラスタインターコネクトを使用するには、Sun Cluster のインストール時に設定したプライベートホスト名をアプリケーションで使用する必要があります。たとえば、host1 のプライベートホスト名が `clusternode1-priv` である場合、クラスタインターコネクトを経由して host1 と通信するときはこの名前を使用する必要があります。この名前を使用してオープンされた TCP ソケットは、クラスタインターコネクトを経由するように経路指定され、プライベートネットワークアダプタに障害が発生した場合でも、この TCP ソケットは透過的に再経路指定されます。任意の 2 つのホスト間のアプリケーション通信はすべてのインターコネクト経由で転送されます。ある特定の TCP 接続のトラフィックは、ある特定の時点では、1 つのインターコネクト上を流れます。異なる TCP 接続はすべてのインターコネクト間で転送されます。さらに、UDP トラフィックは常に、すべてのインターコネクト間で利用されます。

アプリケーションでは、ゾーンのプライベートホスト名を使用し、ゾーン間でクラスタインターコネクトを経由して通信することもできます。ただし、各ゾーンのプライベートホスト名をまず設定しないと、アプリケーションは通信を開始することができません。各ゾーンは、通信するための独自のプライベートホスト名を持っている必要があります。あるゾーンで動作中のアプリケーションは、同じゾーン内のプライベートホスト名を使用して、他ゾーン内のプライベートホスト名と通信します。1 つのゾーン内の各アプリケーションは、別のゾーン内のプライベートホスト名を使って通信することはできません。

Sun Cluster のインストール時に複数のプライベートホスト名が設定されているため、クラスタインターコネクトでは、そのときに選択した任意の名前を使用できます。実際の名前を判別するために `scha_cluster_get` コマンドを `scha_privatelink_hostname_node` 引数と組み合わせて使用します。詳細は、[scha\\_cluster\\_get\(1HA\)](#) のマニュアルページを参照してください。

各ホストには、固定した `per-host` アドレスが割り当てられます。この `per-host` アドレスは、`clprivnet` ドライバで探索されます。IP アドレスは、ホストのプライベートホスト名にマッピングされます。つまり、`clusternode1-priv` です。詳細は、[clprivnet\(7\)](#) のマニュアルページを参照してください。

アプリケーション全体で一貫した IP アドレスが必要な場合、クライアント側とサーバー側の両方でこの `per-host` アドレスにバインドするようにアプリケーションを設定します。これによって、すべての接続はこの `per-host` アドレスから始まり、そして戻されるように見えます。

## リソース、リソースグループ、リソースタイプ

データサービスは、複数の「リソース」タイプを利用します。Sun Java System Web Server や Apache Web Server などのアプリケーションは、それらのアプリケーションが依存するネットワークアドレス (論理ホスト名と共有アドレス) を使用します。アプリケーションとネットワークリソースは RGM が管理する基本単位です。

データサービスはリソースタイプです。たとえば、Sun Cluster HA for Oracle のリソースタイプは SUNW.oracle-server、Sun Cluster HA for Apache のリソースタイプ SUNW.apache です。

リソースはリソースタイプをインスタンス化したもので、クラスタ規模で定義されます。いくつかのリソースタイプはすでに定義されています。

ネットワークリソースは、SUNW.LogicalHostname リソースタイプまたは SUNW.SharedAddress リソースタイプのどちらかです。これら 2 つのリソースタイプは、Sun Cluster ソフトウェアにより事前に登録されています。

HASStoragePlus リソースタイプは、リソースとそのリソースが依存するディスクデバイスグループの起動を同期させるのにも使用できます。このリソースタイプによって、クラスタファイルシステムのマウントポイント、広域デバイス、およびデバイスグループ名のパスがデータサービスの起動前に利用可能になることが保証されます。詳細は、『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』の「リソースグループとデバイスグループ間での起動の同期」を参照してください。

。HASStoragePlus リソースタイプはまた、ローカルのファイルシステムを高可用対応にすることができます。この機能についての詳細は、58 ページの「HASStoragePlus リソースタイプ」を参照してください。

RGM が管理するリソースは、1 つのユニットとして管理できるように、「リソースグループ」と呼ばれるグループに配置されます。リソースグループ上でフェイルオーバーまたはスイッチオーバーが開始されると、リソースグループは 1 つのユニットとして移行されます。

---

注-アプリケーションリソースが含まれるリソースグループをオンラインにすると、そのアプリケーションが起動します。データサービスの起動メソッドは、アプリケーションが起動され、実行されるのを待ってから、正常に終了します。アプリケーションの起動と実行のタイミングの確認は、データサービスがクライアントにサービスを提供しているかどうかをデータサービスの障害モニターが確認する方法と同じです。このプロセスについての詳細は、『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』を参照してください。

---

## リソースグループマネージャー (Resource Group Manager、RGM)

RGMは、データサービス(アプリケーション)を、リソースタイプの実装によって管理されるリソースとして制御します。これらの実装は、Sunが行う場合もあれば、開発者が汎用データサービステンプレート、データサービス開発ライブラリ API (DSDL API)、またはリソース管理 API (RMAPI) を使用して作成することもあります。クラスタ管理者は、「リソースグループ」と呼ばれる入れ物(コンテナ)の中でリソースの作成や管理を行います。RGMは、クラスタメンバーシップの変更に応じて、指定ノードのリソースグループを停止および開始します。

RGMは「リソース」と「リソースグループ」に作用します。RGMのアクションによって、リソースやリソースグループの状態はオンラインまたはオフラインに切り替えられます。リソースとリソースグループに適用できる状態と設定値についての詳細は、[85 ページの「リソースおよびリソースグループの状態と設定値」](#)を参照してください。

RGM 制御下で Solaris プロジェクトを起動する方法については、[95 ページの「データサービスプロジェクトの構成」](#)を参照してください。

## リソースおよびリソースグループの状態と設定値

リソースやリソースグループの値はシステム管理者によって静的に設定されるため、これらの設定は管理操作によってのみ変更できます。RGMは、動的な「状態」の間でリソースグループを移動させます。

これらの設定および状態は次のとおりです。

- **managed** (管理) または **unmanaged** (非管理) 設定。クラスタ全体に適用されるこの設定値は、リソースグループだけに適用されます。リソースグループの管理は RGM が行います。clresourcegroup コマンドを使用して、RGM でリソースグループを管理または非管理するように要求できます。これらのリソースグループ設定は、クラスタ再構成では変更されません。

はじめて作成したリソースグループの状態は非管理になっています。このグループのいずれかのリソースをアクティブにするには、リソースグループの状態が管理になっている必要があります。

スケーラブル Web サーバーなど、ある種のデータサービスでは、ネットワークリソースの起動前や停止後に、あるアクションを行う必要があります。このアクションには、initialization (INIT) と finish (FINI) データサービスメソッドを使用します。INIT メソッドが動作するためには、リソースが置かれているリソースグループが管理状態になっていなければなりません。

リソースグループを非管理から管理の状態に変更すると、そのグループに対して登録されている INIT メソッドがグループの各リソースに対して実行されます。

リソースグループを管理から非管理の状態に変更すると、登録されている FINI メソッドが呼び出され、クリーンアップが行われます。

一般的に、INIT メソッドおよび FINI メソッドは、スケーラブルサービスのネットワークリソース用です。しかし、データサービス開発者は、これらのメソッドをアプリケーションが実行しない初期設定やクリーンアップにも使用できます。

- **enabled** (有効) または **disabled** (無効) 設定。これらの設定は、1 つまたは複数のノード上のリソースに適用されます。システム管理者は、clresource コマンドを使用して、1 つまたは複数のノード上のリソースを有効または無効にできます。これらの設定は、クラスタ管理者がクラスタを再構成しても変わりません。

リソースの通常の設定では、リソースは有効にされ、システムでアクティブに動作しています。

すべてのクラスタノード上でリソースを使用不能にする場合は、すべてのクラスタノード上でリソースを無効にします。無効にしたリソースは、指定したクラスタノード上では、一般的な使用には提供されません。

- **online** (オンライン) または **offline** (オフライン) 状態。動的に変更可能なこれらの状態は、リソースとリソースグループに適用されます。

オンラインとオフラインの状態は、スイッチオーバーまたはフェイルオーバー中、クラスタ再構成手順に従ったクラスタの遷移とともに変化します。さらに clresource および clresourcegroup コマンドを使用して、リソースまたはリソースグループのオンラインまたはオフライン状態を変更することもできます。

フェイルオーバーリソースまたはリソースグループは、常に1つのノード上でのみオンラインにすることができます。スケーラブルリソースまたはリソースグループは、いくつかのノードではオンラインにし、ほかのノードではオフラインにすることができます。スイッチオーバーまたはフェイルオーバー時には、リソー

スグループとそれに含まれるリソースは、あるノードでオフラインになり、その後、別のノードでオンラインになります。

リソースグループがオフラインの場合、そのリソースグループのすべてのリソースがオフラインです。リソースグループがオンラインの場合、そのリソースグループのすべてのリソースがオンラインです。

リソースグループの自動回復アクションを一時的に中断することもできます。リソースグループの自動復旧は、クラスタ内にある問題を調査して修正するために、中断する必要がある場合があります。または、リソースグループサービスの保守作業を実行しなければならない場合もあります。

自動復旧を再開するコマンドを明示的に実行するまで、中断されたリソースグループが自動的に再開またはフェイルオーバーされることはありません。中断されたデータサービスは、オンラインかオフラインかにかかわらず、現在の状態のままとなります。指定されたノード上では、この状態でもリソースグループを別の状態に手作業で切り替えられます。また、リソースグループ内の個々のリソースも有効または無効にできます。

リソースグループはいくつかのリソースを持つことができますが、リソース間には相互依存関係があります。したがって、これらのリソースをオンラインまたはオフラインにするときには、特定の順序で行う必要があります。リソースをオンラインまたはオフラインにするためにメソッドが必要とする時間は、リソースによって異なります。リソースの相互依存関係と起動や停止時間の違いにより、クラスタの再構成では、同じリソースグループのリソースでもオンラインやオフラインの状態が異なる場合があります。

## リソースとリソースグループプロパティ

Sun Cluster データサービスのリソースやリソースグループのプロパティ値は構成できます。標準的なプロパティはすべてのデータサービスに共通です。拡張プロパティは各データサービスに特定のもので、標準プロパティおよび拡張プロパティのいくつかは、デフォルト設定によって構成されているため、これらを修正する必要はありません。それ以外のプロパティは、リソースを作成して構成するプロセスの一部として設定する必要があります。各データサービスのマニュアルでは、設定できるリソースプロパティの種類とその設定方法を指定しています。

標準プロパティは、通常特定のデータサービスに依存しないリソースおよびリソースグループプロパティを構成するために使用されます。標準プロパティのセットについては、『[Sun Cluster データサービスの計画と管理 \(Solaris OS 版\)](#)』の付録 B「標準プロパティ」を参照してください。

RGM 拡張プロパティは、アプリケーションバイナリの場所や構成ファイルなどの情報を提供するものです。拡張プロパティは、データサービスの構成に従って修正する必要があります。拡張プロパティについては、データサービスの個別のガイドで説明されています。

## Solaris ゾーンをサポート

Solaris ゾーンは、Solaris 10 OS のインスタンス内で仮想化されたオペレーティングシステム環境を作成する手段を提供します。Solaris ゾーンを使用すると、1つまたは複数のアプリケーションをシステム上のほかの活動から分離して実行できます。Solaris ゾーンの機能については、『[Solaris のシステム管理 \(Solaris コンテナ: 資源管理と Solaris ゾーン\)](#)』のパート II 「ゾーン」を参照してください。

Solaris 10 OS 上で Sun Cluster ソフトウェアを実行する場合は、任意の数のグローバルクラスタ非投票ノードを作成できます。

Sun Cluster ソフトウェアを使用して、グローバルクラスタ非投票ノード上で動作するアプリケーションの可用性とスケラビリティを管理できます。

## RGM によるグローバルクラスタ非投票ノード (Solaris ゾーン) の直接サポート

Solaris 10 OS が動作するクラスタでは、リソースグループをグローバルクラスタ投票ノードまたはグローバルクラスタ非投票ノードで動作するように設定できます。RGM は、各グローバルクラスタ非投票ノードをスイッチオーバーターゲットとして管理します。リソースグループのノードリストでグローバルクラスタ非投票ノードが指定されている場合、RGM は指定されたノードでリソースグループをオンラインにします。

図 3-8 は、2 ホストクラスタでのノード間のリソースグループのフェイルオーバーを示しています。この例では、クラスタの管理を簡単にするために同一のノードが構成されています。



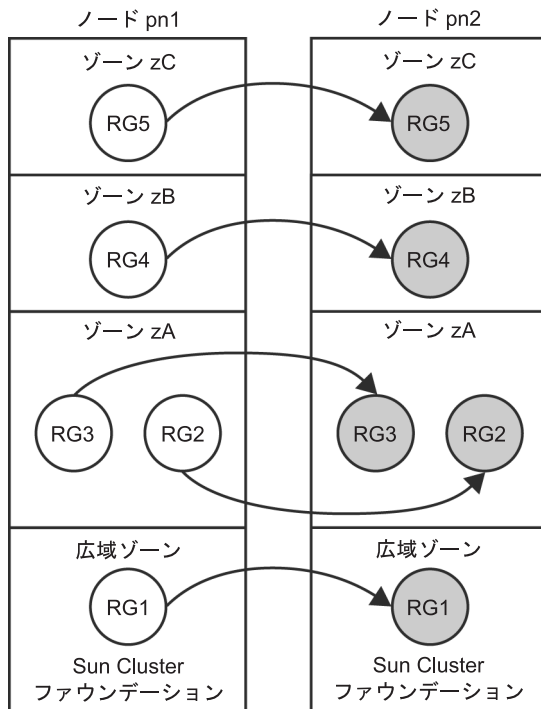


図 3-8 ノード間のリソースグループのフェイルオーバー

スケーラブルなリソースグループ(ネットワーク負荷分散を使用)を、クラスタ非投票ノードでも動作するよう構成することができます。

Sun Cluster コマンドで、次の例に示すように、ゾーンの名前をホストの名前に追加し、それらをコロンで区切ることによって、ゾーンを指定します。

```
phys-schost-1:zoneA
```

## RGM による Solaris ゾーン の直接サポート を使用する ための 基準

次のいずれかの基準を満たす場合、RGM による Solaris ゾーン の直接サポート を使用 します。

- アプリケーションがゾーンを起動するために必要な追加のフェイルオーバー時間を許容できない。
- メンテナンス中に停止時間を最小にする必要がある。
- デュアルパーティションのソフトウェアアップグレードを必要としている。
- ネットワーク負荷分散のために共有アドレスリソースを使用するデータサービスを構成している。

## RGM による Solaris ゾーン の直接サポート を使用するための要件

アプリケーションで RGM による Solaris ゾーン の直接サポート の使用を計画している場合は、次の要件を満たしていることを確認してください。

- アプリケーションが非大域ゾーンでの動作に対応していること。
- アプリケーションのデータサービスがグローバルクラスタ非投票ノードでの動作に対応していること。

RGM による Solaris ゾーン の直接サポート を使用する場合、アフィニティーにより関連付けられているリソースグループが同じ Solaris ホストで動作するように構成されていることを確認します。

## RGM による Solaris ゾーン の直接サポート に関するその他の情報

RGM による Solaris ゾーン の直接サポート の構成方法の詳細については、次のドキュメントを参照してください。

- 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』の「グローバルクラスタ内の非大域ゾーンのガイドライン」
- 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』の「ゾーン名」
- 『Sun Cluster ソフトウェアのインストール (Solaris OS 版)』の「グローバルクラスタノード上での非大域ゾーンの設定」
- 『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』
- 各データサービスガイド

## Sun Cluster ノード上の Solaris ゾーン を Sun Cluster HA for Solaris Containers を通してサポート

Sun Cluster HA for Solaris Containers データサービスは、各ゾーンを RGM によって制御されるリソースとして管理します。

## Sun Cluster HA for Solaris Containers を使用するための基準

次のいずれかの基準を満たす場合、Sun Cluster HA for Solaris Containers データサービスを使用します。

- 代理のルートアクセスを必要とする。
- アプリケーションがクラスタでサポートされていない。
- 同じノードの別のゾーンで動作するリソースグループ間にアフィニティーを必要としている。

## Sun Cluster HA for Solaris Containers を使用するための要件

アプリケーションで Sun Cluster HA for Solaris Containers データサービスの利用を計画している場合、次の要件を満たすことを確認してください。

- アプリケーションがグローバルクラスタ非投票ノードでの動作に対応していること。
- アプリケーションがスクリプト、実行レベルのスクリプト、または Solaris サービス管理機能 (Service Management Facility、SMF) のマニフェストによって Solaris OS と統合されていること。
- ゾーンを起動するために必要な追加のフェイルオーバー時間を許容できること。
- メンテナンス中の停止時間を許容できること。

## Sun Cluster HA for Solaris Containers についてのその他の情報

Sun Cluster HA for Solaris Containers データサービスの使い方の詳細については、『[Sun Cluster Data Service for Solaris Containers Guide for Solaris OS](#)』を参照してください。

# サービス管理機能

Solaris サービス管理機能 (Service Management Facility、SMF) によって、アプリケーションを可用性が高く、スケーラブルなリソースとして実行して管理することができます。リソースグループマネージャー (Resource Group Manager、RGM) と同じように、SMF は高可用性とスケーラビリティを提供しますが、Solaris オペレーティングシステム用です。

Sun Cluster は、クラスタで SMF サービスを有効にするために使用する 3 種類のプロキシリソースタイプを提供します。これらのリソースタイプ、`SUNW.Proxy_SMF_failover`、`SUNW.Proxy_SMF_loadbalanced`、および `SUNW.Proxy_SMF_multimaster` により、フェイルオーバー、スケーラブル、およびマルチマスターのそれぞれの構成で、SMF サービスを実行できます。SMF は単一 Solaris ホスト上で SMF サービスの可用性を管理します。SMF はコールバックメソッド実行モデルを使用して、サービスを実行します。

さらに、SMF はサービスの監視と制御のための管理インタフェースのセットも提供します。これらのインタフェースにより、ユーザー独自の SMF 制御サービスを Sun Cluster に組み込むことができます。この機能により新たなコールバックメソッドを作成したり、既存のコールバックメソッドを書き直したり、あるいは SMF サービスマニフェストを更新する必要がなくなります。複数の SMF リソースを 1 つのリソースグループに含め、それらのリソース間に依存性とアフィニティを構成することができます。

SMF はこれらのサービスを開始、停止、および再開する権限を持ち、サービス間の依存性を管理します。Sun Cluster は、クラスタ内でこれらのサービスを管理し、これらのサービスを開始するホストを決める権限を持ちます。

SMF は、各クラスタホスト上でデーモン `svc.startd` として動作します。SMF デーモンは、あらかじめ設定されたポリシーに基づいて、選択したホスト上で自動的にリソースを開始および停止します。

SMF プロキシリソースに指定されたサービスは、グローバルクラスタ投票ノードまたはグローバルクラスタ非投票ノードに常駐できます。ただし、同じ SMF プロキシリソースに指定したサービスは、同じノードに置く必要があります。SMF プロキシリソースはどのノードでも動作します。

## システムリソースの使用状況

システムリソースには、CPU 使用率、メモリ使用率、スワップ使用率、ディスクおよびネットワークのスループットが含まれます。Sun Cluster では、オブジェクトタイプ別にシステムリソースの使用率を監視できます。オブジェクトタイプには、ホスト、ノード、ゾーン、ディスク、ネットワークインタフェース、またはリソースグループがあります。Sun Cluster ではまた、リソースグループで使用できる CPU を制御することもできます。

システムリソース使用量の監視と制御をリソース管理ポリシーの一部にすることができます。多数のマシンの管理は複雑でコストがかかるため、より大規模なホストにアプリケーションを統合することが望まれます。個々の作業負荷を別々のシステムで実行して、そのシステムのリソースへのフルアクセスを与える代わりに、リソース管理ソフトウェアを使用すれば、システム内の作業負荷を分離できます。リソース管理機能を使用すると、1つの Solaris システムで複数の異なるアプリケーションを実行して制御することにより、システムの総保有コスト (TCO) を低減することができます。

リソース管理機能を使用して、アプリケーションが必要な応答時間を確保できるようにします。また、リソース管理機能により、リソースの使用率を向上させることができます。使用状況を分類して優先付けすることにより、オフピーク時に余った資源を効率よく使用でき、処理能力を追加する必要がなくなります。また、負荷の変動が原因で資源を無駄にすることもなくなります。

Sun Cluster がシステムリソースの使用率について収集するデータを活用するには、次の手順を実行する必要があります。

- データを分析して、システムへの影響を判断する。
- ハードウェアリソースおよびソフトウェアリソースの使用率を最適化するために必要な操作を決定する。
- 決定した操作を実行する。

Sun Cluster のインストール時にデフォルトでは、システムリソースの監視と制御は構成されていません。これらのサービスの構成の詳細は、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の第9章「CPU 使用率の制御の構成」を参照してください。

## システムリソース監視

システムリソースの使用率を監視することにより、次のことができます。

- 特定のシステムリソースを使用するサービスの実行状況を反映するデータを収集する。
- リソースの障害またはオーバーロードを見つけて、事前に問題を回避する予防策を取る。
- より効率的にワークロードを管理する。

システムリソースの使用率に関するデータにより、あまり使用されていないハードウェアリソースや多くのリソースを使用するアプリケーションを判別することができます。このデータに基づいて、必要なリソースを備えたノードにアプリケーションを割り当て、フェイルオーバーするノードを選択することができます。この統合により、ハードウェアリソースとソフトウェアリソースの使用方法を最適化できます。

すべてのシステムリソースを同時に監視することは、CPU の負担になる場合があります。システムにもっとも重要なリソースに優先順位を付けて、監視するシステムリソースを選択してください。

監視を有効にするときに、監視するテレメトリ属性を選択します。テレメトリ属性はシステムリソースの一面です。テレメトリ属性の例としては、CPU の量またはデバイスで使用されているブロックの使用率などがあります。あるオブジェクトタイプについてテレメトリ属性を監視する場合、Sun Cluster はクラスタ内のそのタイプのすべてのオブジェクトでこのテレメトリ属性を監視します。Sun Cluster は収集されるシステムリソースデータを7日間保存します。

特定のデータ値がシステムリソースに重要だと考えられる場合、この値にしきい値を設定できます。しきい値を設定する際には、しきい値に重要度を割り当てることにより、このしきい値のクリティカル度の選択も行います。このしきい値を超えると、Sun Cluster はこのしきい値の重要度レベルをユーザーが選択する重要度レベルに変更します。

## CPU の制御

クラスタ上で動作するアプリケーションおよびサービスウィンドウごとに特定の CPU ニーズがあります。表 3-4 は、Solaris OS の各バージョンで使用できる CPU 制御動作を示しています。

表 3-4 CPU の制御

| Solaris のバージョン | ゾーン             | 制御  |
|----------------|-----------------|---|
| Solaris 9 OS   | 使用不可            | CPU の配分の割り当て                                  |
| Solaris 10 OS  | グローバルクラスタ投票ノード  | CPU の配分の割り当て                                  |
| Solaris 10 OS  | グローバルクラスタ非投票ノード | CPU の配分の割り当て<br>CPU の数の割り当て<br>専用のプロセッサセットの作成 |

注 - CPU シェアを提供する場合、クラスタ内でフェアシェアスケジューラ (FFS) をデフォルトのスケジューラとして指定する必要があります。

グローバルクラスタ非投票ノードで専用プロセッサセットのリソースグループに割り当てられている CPU を制御することにより、もっとも厳格なレベルの制御が実現されます。あるリソースグループに CPU を予約すると、この CPU はほかのリソースグループでは使用できません。

## システムリソース使用率の表示

コマンドラインまたは Sun Cluster Manager を使用して、システムリソースデータおよび CPU 割り当てを表示できます。監視を選択するシステムリソースによって、表示できる表とグラフが決まります。

システムリソース使用率と CPU 制御の出力を表示することにより、次を実現することができます。

- システムリソースの消耗による障害を予測する。
- システムリソースの使用率の不均衡を見つける。
- サーバーの統合性を確認する。
- アプリケーションのパフォーマンスを改善できる情報を取得する。

Sun Cluster では、収集したデータに基づいて取るべき措置をアドバイスしたり、ユーザーの代わりに措置を実行することはありません。表示されるデータがサービスに期待する条件を満たしているかどうかを判断する必要があります。そのあとで、確認されたパフォーマンスの救済措置を取る必要があります。

# データサービスプロジェクトの構成

データサービスは、RGMでオンラインにしたときに Solaris プロジェクト名のもとで起動するように構成できます。そのためには、データサービスを構成するときに、RGMによって管理されるリソースまたはリソースグループと Solaris プロジェクト ID を対応付ける必要があります。リソースまたはリソースグループにプロジェクト ID を対応付けることによって、Solaris オペレーティングシステムの洗練されたコントロールを使用して、クラスタ内の負荷や使用量を管理できます。

---

注 - Solaris 9 OS または Solaris 10 OS 上で Sun Cluster を使用する場合、この構成を実行できます。

---

Sun Cluster 環境の Solaris 管理機能を使用すると、ほかのアプリケーションとノードを共有している場合に、もっとも重要なアプリケーションに高い優先順位を与えることができます。ノードを複数のアプリケーションで共有する例としては、サービスを統合した場合や、アプリケーションのフェイルオーバーが起きた場合があります。ここで述べる管理機能を使用すれば、優先順位の低いアプリケーションが CPU 時間などのシステムサプライを過度に使用するのを防止し、重要なアプリケーションの可用性を向上させることができます。

---

注 - この機能に関連する Solaris のマニュアルでは、CPU 時間、プロセス、タスクなどのコンポーネントを「リソース」と呼んでいます。一方、Sun Cluster のマニュアルでは、RGM の制御下にあるエンティティを「リソース」と呼んでいます。次の節では、「リソース」という用語を RGM で制御される Sun Cluster エンティティを指す用語として使用します。また、CPU 時間、プロセス、およびタスクを「サプライ」と呼びます。

---

この節では、指定した Solaris OS の [project\(4\)](#) でプロセスを起動するように、データサービスを構成する方法の概念について説明します。また、Solaris オペレーティングシステムの管理機能を使用するための、フェイルオーバーのシナリオやヒントについても説明します。

管理機能の概念や手順についての詳細は、『[Solaris のシステム管理\(ネットワークサービス\)](#)』の第 1 章「[ネットワークサービス\(概要\)](#)」を参照してください。

クラスタ内で Solaris 管理機能を使用できるようにリソースやリソースグループを構成するための手順は次のようになります。

1. アプリケーションをリソースの一部として構成します。
2. リソースをリソースグループの一部として構成します。
3. リソースグループのリソースを有効にします。
4. リソースグループを管理状態にします。

5. リソースグループに対する Solaris プロジェクトを作成します。
6. 手順5で作成したプロジェクトにリソースグループ名を対応付けるための標準プロパティを構成します。
7. リソースグループをオンラインにします。

Solaris プロジェクト ID をリソースやリソースグループに関連付けるように標準の `Resource_project_name` または `RG_project_name` プロパティを構成するには、`-p` オプションを `clresource set` および `clresourcegroup set` コマンドとともに使用します。続いて、プロパティの値にリソースまたはリソースグループを設定します。プロパティの定義については、『[Sun Cluster データサービスの計画と管理 \(Solaris OS 版\)](#)』の付録 B 「標準プロパティ」を参照してください。プロパティの説明については、`r_properties(5)` および `rg_properties(5)` のマニュアルページを参照してください。

指定するプロジェクト名はプロジェクトデータベース (`/etc/project`) に存在するものでなければなりません。さらに、指定するプロジェクトのメンバーとして `root` ユーザーが設定されていなければなりません。プロジェクト名データベースの概念については、『[Solaris のシステム管理 \(Solaris コンテナ: 資源管理と Solaris ゾーン\)](#)』の第2章「プロジェクトとタスク (概要)」を参照してください。プロジェクトファイルの構文については、`project(4)` を参照してください。

RGM は、リソースまたはリソースグループをオンラインにする際に、関連するプロセスをこのプロジェクト名の下で起動します。

---

注-リソースまたはリソースグループとプロジェクトを対応付けることはいつでもできます。ただし、RGM を使ってプロジェクトのリソースやリソースグループをオフラインにしてから再びオンラインに戻すまで、新しいプロジェクト名は有効になりません。

---

リソースやリソースグループをプロジェクト名の下で起動すれば、次の機能を構成することによってクラスタ全体のシステムサプライを管理できます。

- 拡張アカウンティング-使用量をタスクやプロセス単位で記録できるため柔軟性が増します。拡張アカウンティングでは、使用状況の履歴を調べ、将来の作業負荷の容量要件を算定できます。
- 制御-システムサプライの使用を制約する機構を提供します。これにより、プロセス、タスク、およびプロジェクトが特定のシステムサプライを大量に消費することを防止できます。
- フェアシェアスケジューリング (FSS)-それぞれの作業負荷に割り当てる CPU 時間を作業負荷の重要性に基づいて制御できます。作業負荷の重要性は、各作業負荷に割り当てる、CPU 時間のシェア数として表されます。詳細は、次のマニュアルページを参照してください。

- `dispadmin(1M)`



- `priocntl(1)`
- `ps(1)`
- `FSS(7)`
- プール-アプリケーションの必要性に応じて対話型アプリケーション用に仕切りを使用することができます。プールを使用すれば、ホストを仕切り分けすることができます、同じサーバーで異なるソフトウェアアプリケーションをサポートできます。プールを使用すると、アプリケーションごとの応答が予測しやすくなります。

## プロジェクト構成に応じた要件の決定

Sun Cluster 環境で Solaris が提供する制御を使用するようにデータサービスを構成するには、まず、スイッチオーバーやフェイルオーバー時にリソースをどのように制御および管理するかを決めておく必要があります。新しいプロジェクトを構成する前に、まず、クラスタ内の依存関係を明確にします。たとえば、リソースやリソースグループはデバイスグループに依存しています。

リソースグループのノードリストの優先順位を識別するには、`nodelist`、`failback`、`maximum primaries` および `desired primaries` リソースグループのプロパティを使用します。これらのプロパティは、`clresourcegroup set` コマンドで構成します。

- リソースグループとデバイスグループのノードリストの依存性に関する簡単な説明については、『[Sun Cluster データサービスの計画と管理 \(Solaris OS 版\)](#)』の「[リソースグループとデバイスグループの関係](#)」を参照してください。
- プロパティの詳細な説明については、[rg\\_properties\(5\)](#) を参照してください。

デバイスグループのノードリストの優先順位を決めるには、`preferenced` プロパティおよび `failback` プロパティを使用します。これらのプロパティは、`cldevicegroup` および `clsetup` コマンドで構成します。詳細は、[clresourcegroup\(1CL\)](#)、[cldevicegroup\(1CL\)](#)、および [clsetup\(1CL\)](#) のマニュアルページを参照してください。

- `preferenced` プロパティの概念については、53 ページの「[多重ポートデバイスグループ](#)」を参照してください。
- 手順については、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の「[デバイスグループの管理](#)」の「[ディスクデバイスのプロパティを変更する](#)」を参照してください。
- ノード構成の概念とフェイルオーバーデータサービスとスケラブルデータサービスの動作については、21 ページの「[Sun Cluster システムのハードウェアおよびソフトウェアコンポーネント](#)」を参照してください。

すべてのクラスタノードを同じように構成すると、主ノードと二次ノードに対して同じ使用限度が割り当てられます。各プロジェクトの構成パラメータは、すべてのノードの構成ファイルに定義されているすべてのアプリケーションに対して同じで

ある必要はありません。特定のアプリケーションに対応するすべてのプロジェクトは、少なくとも、そのアプリケーションのすべての潜在的マスターにあるプロジェクトデータベースからアクセス可能である必要があります。アプリケーション1が *phys-schost-1* によってマスターされているが、*phys-schost-2* または *phys-schost-3* に切り替えられるか、フェイルオーバーされる可能性があるかと仮定します。アプリケーション1に対応付けられたプロジェクトは、これら3つのノード (*phys-schost-1*、*phys-schost-2*、*phys-schost-3*) でアクセス可能でなければなりません。

---

注- プロジェクトデータベース情報は、ローカルの `/etc/project` データベースファイルに格納することも、NIS マップや LDAP ディレクトリサーバーに格納することもできます。

---

Solaris オペレーティングシステムでは、使用パラメータは柔軟に構成でき、Sun Cluster によって課せられる制約はほとんどありません。どのような構成を選択するかはサイトの必要性によって異なります。システムの構成を始める前に、次の各項の一般的な指針を参考にしてください。

## プロセス当たりの仮想メモリー制限の設定

仮想メモリーの制限をプロセス単位で制御する場合は、`process.max-address-space` コントロールを使用します。`process.max-address-space` 値の設定方法についての詳細は、`rctldm(1M)` のマニュアルページを参照してください。

Sun Cluster ソフトウェアで管理コントロールを使用する場合は、アプリケーションの不要なフェイルオーバーが発生したり、アプリケーションの「ピンポン」現象が発生したりするのを防止するために、メモリー制限を適切に設定する必要があります。そのためには、一般に次の点に注意する必要があります。

- メモリー制限をあまり低く設定しない。  
アプリケーションは、そのメモリーが限界に達すると、フェイルオーバーを起こすことがあります。データベースアプリケーションにとってこの指針は特に重要です。その仮想メモリーが限界を超えると予期しない結果になることがあるからです。
- 主ノードと二次ノードに同じメモリー制限を設定しない。  
同じメモリー制限を設定すると、アプリケーションのメモリーが限度に達し、アプリケーションが、同じメモリー制限をもつ二次ノードにフェイルオーバーされたときに「ピンポン」現象を引き起こす可能性があります。そのため、二次ノードのメモリー制限には、主ノードよりもわずかに大きな値を設定します。異なるメモリー制限を設定することによって「ピンポン」現象の発生を防ぎ、管理者はその間にパラメータを適切に変更することができます。
- 負荷均衡を達成する目的でリソース管理メモリー制限を使用する。

たとえば、メモリ制限を使用すれば、アプリケーションが誤って過度のスワップ領域を使用することを防止できます。

## フェイルオーバーシナリオ

管理パラメータを適切に構成すれば、プロジェクト構成 (/etc/project) 内の割り当ては、通常のクラスタ操作でも、スイッチオーバーやフェイルオーバーの状況でも正常に機能します。

以下の各項ではシナリオ例を説明します。

- 最初の2つの項、100ページの「2つのアプリケーションを供う2ホストクラスタ」および101ページの「3つのアプリケーションを供う2ホストクラスタ」では、ホスト全体が関係するフェイルオーバーシナリオについて説明します。
- 103ページの「リソースグループだけのフェイルオーバー」の項では、アプリケーションだけのフェイルオーバー操作について説明します。

Sun Cluster 環境では、アプリケーションはリソースの一部として構成します。そして、リソースをリソースグループ (RG) の一部として構成します。障害が発生すると、リソースグループは、対応付けられたアプリケーションとともに、別のノードにフェイルオーバーされます。以下の例では、リソースは明示的に示されていません。各リソースには、1つのアプリケーションが構成されているものとします。

---

注-フェイルオーバーは、ノードがノードリストで指定され、RGMで設定された順序で行なわれます。

---

以下の例は次のように構成されています。

- アプリケーション1 (App-1) はリソースグループ RG-1 に構成されています。
- アプリケーション2 (App-2) はリソースグループ RG-2 に構成されています。
- アプリケーション3 (App-3) はリソースグループ RG-3 に構成されています。

フェイルオーバーが起こると、各アプリケーションに割り当てられる CPU 時間の割合が変化します。ただし、割り当てられているシェアの数はそのままです。この割合は、そのノードで動作しているアプリケーションの数と、アクティブな各アプリケーションに割り当てられているシェアの数によって異なります。

これらのシナリオでは、次のように構成が行われているものとします。

- すべてのアプリケーションが共通のプロジェクトの下に構成されています。
- 各リソースには1つのアプリケーションがあります。
- すべてのノードにおいて、アクティブなプロセスはこれらのアプリケーションだけです。
- プロジェクトデータベースは、クラスタの各ノードで同一に構成されています。

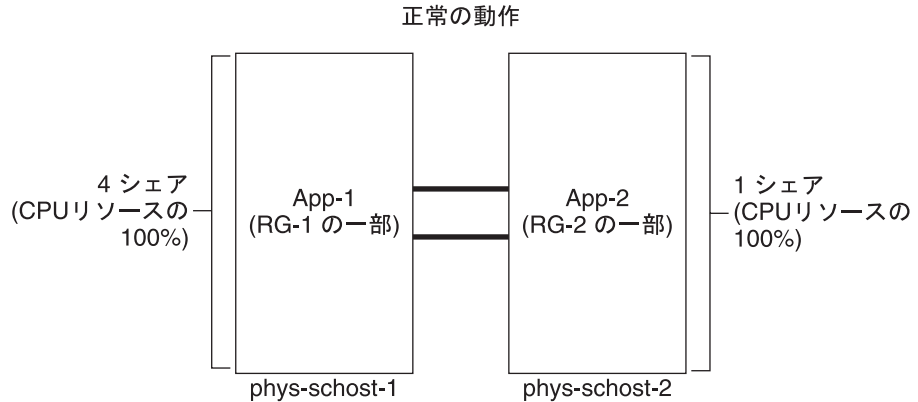
## 2つのアプリケーションを供う2ホストクラスタ

2ホストクラスタに2つのアプリケーションを構成することによって、それぞれの物理ホスト (*phys-schost-1*、*phys-schost-2*) を1つのアプリケーションのデフォルトマスターにすることができます。一方の物理ホストは、他方の物理ホストの二次ノードになります。アプリケーション1とアプリケーション2に関連付けられているすべてのプロジェクトは、両ノードのプロジェクトデータベースファイルに存在している必要があります。クラスタが正常に動作している間、各アプリケーションはそれぞれのデフォルトマスターで動作し、管理機能によってすべてのCPU時間を割り当てられます。

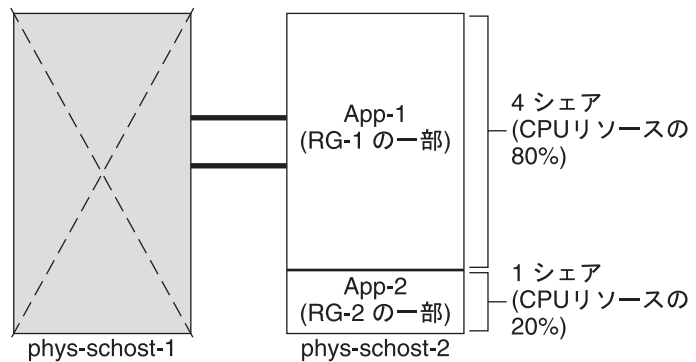
フェイルオーバーかスイッチオーバーが起ると、これらのアプリケーションは同じノードで動作し、構成ファイルの設定に従ってシェアを割り当てられます。たとえば、`/etc/project` ファイルに次のエントリが指定されていると、アプリケーション1に4シェアが、アプリケーション2に1シェアがそれぞれ割り当てられます。

```
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)
```

次の図は、この構成の正常時の動作とフェイルオーバー時の動作を表しています。割り当てられているシェアの数は変わりません。ただし、各アプリケーションが利用できるCPU時間の割合は変わる場合があります。この割合は、CPU時間を要求する各プロセスに割り当てられているシェア数によって異なります。



フェイルオーバー時の動作: ノード phys-schost-1 の障害



### 3つのアプリケーションを供う2ホストクラスタ

3つのアプリケーションが動作する2ホストクラスタでは、1つのホスト (*phys-schost-1*) を1つのアプリケーションのデフォルトマスターとして構成できます。そして、もう1つの物理ホスト (*phys-schost-2*) をほかの2つのアプリケーションのデフォルトマスターとして構成できます。各ホストには、次のサンプルプロジェクトデータベースファイルがあるものとします。フェイルオーバーやスイッチオーバーが起っても、プロジェクトデータベースファイルが変更されることはありません。

```
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)
```

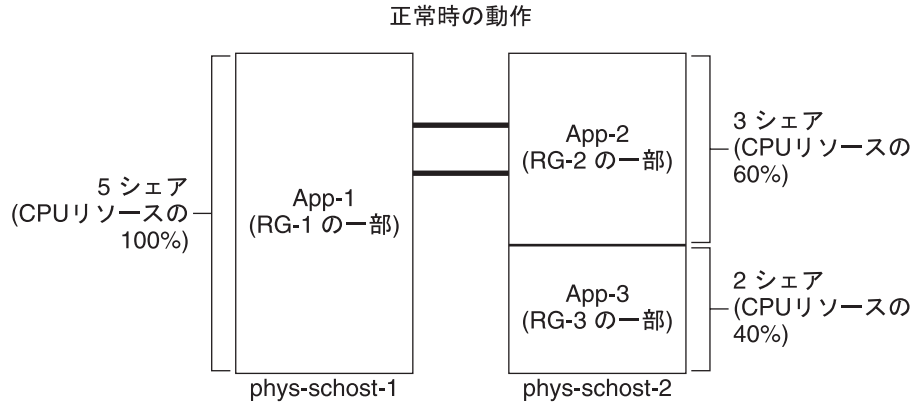
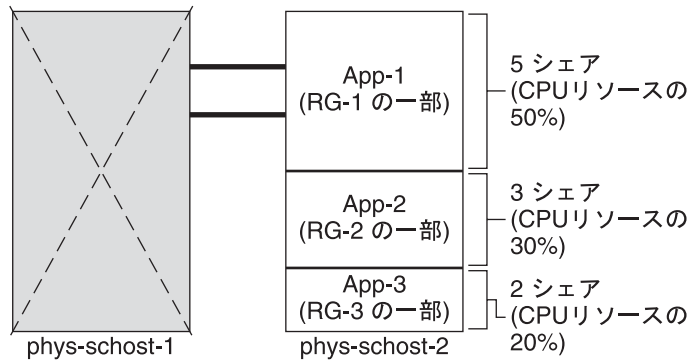
クラスタが正常に動作している間、アプリケーション1には、そのデフォルトマスター *phys-schost-1* で5シェアが割り当てられます。このホストでCPU時間を要求するアプリケーションはこのアプリケーションだけであるため、この数は100パーセントのCPU時間と同じことです。アプリケーション2と3には、それぞれのデフォル

トマスターである *phys-schost-2* で3シェアと2シェアが割り当てられます。したがって、正常な動作では、アプリケーション2にCPU時間の60パーセントが、アプリケーション3にCPU時間の40パーセントがそれぞれ割り当てられます。

フェイルオーバーかスイッチオーバーが発生し、アプリケーション1が *phys-schost-2* に切り替えられても、3つのアプリケーションの各シェアは変わりません。ただし、割り当てられるCPUリソースの割合はプロジェクトデータベースファイルに従って変更されます。

- 5シェアをもつアプリケーション1にはCPUの50パーセントが割り当てられます。
- 3シェアをもつアプリケーション2にはCPUの30パーセントが割り当てられます。
- 2シェアをもつアプリケーション3にはCPUの20パーセントが割り当てられます。

次の図は、この構成の正常な動作とフェイルオーバー動作を示しています。

フェイルオーバー時の動作: ノード `phys-schost-1` の障害

## リソースグループだけのフェイルオーバー

複数のリソースグループが同じデフォルトマスターに属している構成では、1つのリソースグループ(および、それに関連付けられたアプリケーション)が二次ノードにフェイルオーバーされたり、スイッチオーバーされたりすることがあります。その間、クラスタのデフォルトマスターは動作を続けます。

注-フェイルオーバーの際、フェイルオーバーされるアプリケーションには、二次ホスト上の構成ファイルの指定に従ってリソースが割り当てられます。この例の場合、主ホストと二次ホストのプロジェクトデータベースファイルの構成は同じです。

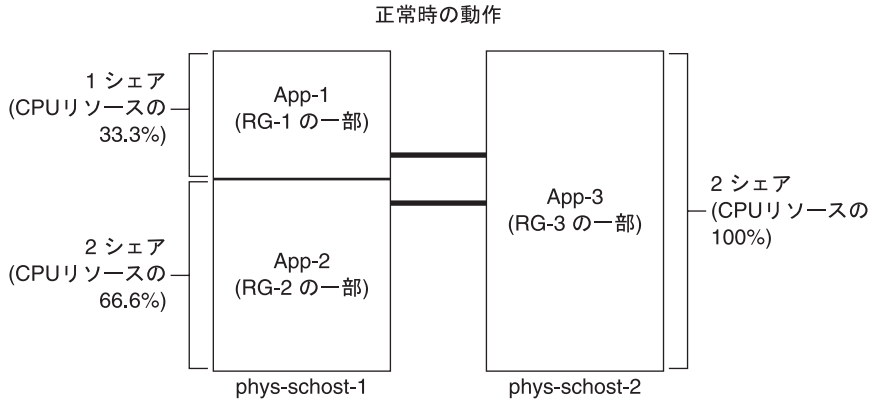
次のサンプル構成ファイルでは、アプリケーション1に1シェア、アプリケーション2に2シェア、アプリケーション3に2シェアがそれぞれ割り当てられています。

```
Prj_1:106:project for App_1:root::project.cpu-shares=(privileged,1,none)
```

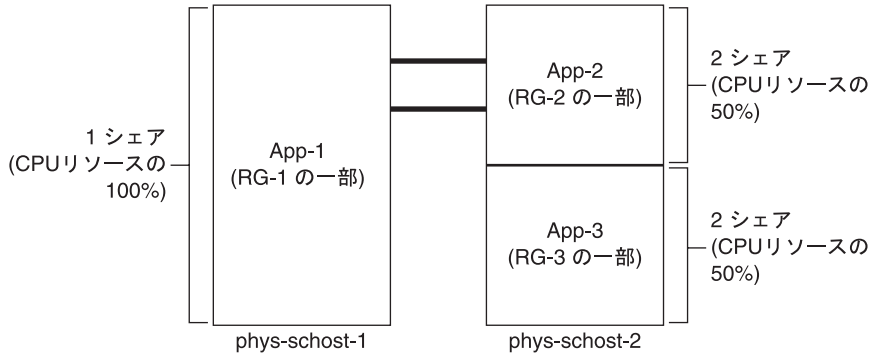
```
Prj_2:107:project for App_2:root::project.cpu-shares=(privileged,2,none)
```

```
Prj_3:108:project for App_3:root::project.cpu-shares=(privileged,2,none)
```

次の図は、この構成の正常時の動作とフェイルオーバー時の動作を表しています。ここでは、アプリケーション2が動作するRG-2が *phys-schost-2* にフェイルオーバーされます。割り当てられているシェアの数は変わりません。ただし、各アプリケーションが利用できるCPU時間の割合は、CPU時間を要求する各アプリケーションに割り当てられているシェア数によって異なります。



フェイルオーバー時の動作: RG-2 の *phys-schost-2* へのフェイルオーバー





# パブリックネットワークアダプタとIPネットワークマルチパス

クライアントは、パブリックネットワークを介してクラスタにデータ要求を行います。各クラスタの Solaris ホストは、1 対のパブリックネットワークアダプタを介して少なくとも1つのパブリックネットワークに接続されています。

Sun Cluster で動作する Solaris インターネットプロトコル (Internet Protocol, IP) ソフトウェアは、パブリックネットワークアダプタを監視したり、障害を検出したときに IP アドレスをあるアダプタから別のアダプタにフェイルオーバーしたりする基本的な機構を提供します。各ホストは独自の IP ネットワークマルチパス 構成を持っており、この構成がほかのホストの構成と異なる場合があります。

パブリックネットワークアダプタは、IP マルチパスグループ(「マルチパスグループ」)として編成されます。各マルチパスグループには、1つまたは複数のパブリックネットワークアダプタがあります。マルチパスグループの各アダプタはアクティブにしておいてもかまいません。あるいは、スタンバイインタフェースを構成し、フェイルオーバーが起こるまでそれらを非アクティブにしておいてもかまいません。

in.mpathd マルチパスデーモンは、テスト IP アドレスを使って障害や修復を検出します。マルチパスデーモンによってアダプタの1つに障害が発生したことが検出されると、フェイルオーバーが行われます。すべてのネットワークアクセスは、障害が発生したアダプタから、マルチパスグループ内の別の動作中のアダプタにフェイルオーバーされます。したがって、デーモンがそのホストのパブリックネットワーク接続を維持します。スタンバイインタフェースを構成していた場合、このデーモンはスタンバイインタフェースを選択します。そうでない場合、このデーモンはもっとも小さい IP アドレス番号を持つインタフェースを選択します。フェイルオーバーはアダプタインタフェースレベルで発生するため、これよりも高いレベルの接続 (TCP など) は影響を受けません。ただし、フェイルオーバー中には一時的にわずかな遅延が発生します。IP アドレスのフェイルオーバーが正常に終了すると、ARP ブロードキャストが送信されます。したがって、デーモンがリモートクライアントへの接続を維持します。

---

注-TCP の輻輳回復特性のために、正常なフェイルオーバーのあと、TCP エンドポイントではさらに遅延が生じる可能性があります。これは、フェイルオーバー中にいくつかのセグメントが失われて、TCP の輻輳制御機構がアクティブになるためです。

---

マルチパスグループには、論理ホスト名と共有アドレスリソースの構築ブロックがあります。論理ホスト名と共有アドレスリソースとは別にマルチパスグループを作成して、クラスタホストのパブリックネットワーク接続を監視する必要もあります。つまり、ホスト上の同じマルチパスグループは、任意の数の論理ホスト名または

共有アドレスリソースをホストできます。論理ホスト名と共有アドレスリソースについての詳細は、『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』を参照してください。

注-IP ネットワークマルチパス 機構の設計は、アダプタの障害を検出してマスクすることを目的としています。この設計は、管理者が `ifconfig` を使用して論理(または共有)IP アドレスのどれかを削除した状態から回復することを目的としているわけではありません。Sun Cluster ソフトウェアから見ると、論理アドレスや共有 IP アドレスは RGM によって管理されるリソースです。管理者が IP アドレスを追加または削除する場合、正しくは、`clresource` および `clresourcegroup` を使用して、リソースを含むリソースグループを修正します。

IP ネットワークマルチパスの Solaris の実装についての詳細は、クラスタにインストールされている Solaris オペレーティングシステムのマニュアルを参照してください。

| オペレーティングシステム            | 参照先   |
|-------------------------|---|
| Solaris 9 オペレーティングシステム  | 『IP ネットワークマルチパスの管理』の第 1 章「IP ネットワークマルチパス(概要)」 |
| Solaris 10 オペレーティングシステム | 『Solaris のシステム管理 (IP サービス)』のパート VI 「IPMP」     |

## SPARC: 動的再構成のサポート

Sun Cluster 3.2 1/09 による動的再構成 (Dynamic Reconfiguration、DR) ソフトウェア機能のサポートは段階的に開発されています。この節では、Sun Cluster 3.2 1/09 による DR 機能のサポートの概念と考慮事項について説明します。

Solaris の DR 機能の説明で述べられているすべての必要条件、手順、制限は Sun Cluster の DR サポートにも適用されます (オペレーティング環境での休止状態中を除く)。したがって、Sun Cluster ソフトウェアで DR 機能を使用する前には、必ず、Solaris の DR 機能についての説明を参照してください。特に、DR Detach 操作中に、ネットワークに接続されていない入出力デバイスに影響する問題について確認してください。

『Sun Enterprise 10000 Dynamic Reconfiguration ユーザーマニュアル』と『Sun Enterprise 10000 Dynamic Reconfiguration リファレンスマニュアル』が、<http://docs.sun.com> で参照できます。

## SPARC: 動的再構成の概要

DR 機能を使用すると、システムハードウェアの切り離しなどの操作をシステムの稼動中に行うことができます。DR プロセスの目的は、システムを停止したり、クラスタの可用性を中断したりせずにシステム操作を継続できるようにすることです。

DR はボードレベルで機能します。したがって、DR 操作はボード上のすべてのコンポーネントに影響します。ボードには、CPU やメモリー、ディスクドライブやテープドライブ、ネットワーク接続の周辺機器インタフェースなど、複数のコンポーネントが取り付けられています。

アクティブなコンポーネントを含むボードを切り離すと、システムエラーになります。DR サブシステムは、ボードを切り離す前に、ほかのサブシステム (Sun Cluster など) に問い合わせたボード上のコンポーネントが使用されているかを判別します。ボードが使用中であることがわかると、DR のボード切り離し操作は行われません。つまり、アクティブなコンポーネントを含むボードに DR のボード切り離し操作を発行しても、DR サブシステムがその操作を拒否するため、DR のボード切り離し操作はいつ発行しても安全です。

同様に、DR のボード追加操作も常に安全です。新たに追加されたボードの CPU とメモリーは、システムによって自動的にサービス状態になります。ただし、そのボードのほかのコンポーネントを意図的に使用するには、管理者がそのクラスタを手動で構成する必要があります。

---

注 - DR サブシステムにはいくつかのレベルがあります。下位のレベルがエラーを報告すると、上位のレベルもエラーを報告します。ただし、下位のレベルが具体的なエラーを報告しても、上位のレベルは「Unknown error」を報告します。このエラーは無視してもかまいません。

---

次の各項では、デバイスタイプごとに DR の注意事項を説明します。

## SPARC: CPU デバイスに対する DR クラスタリング

CPU デバイスが存在していても、Sun Cluster ソフトウェアは DR のボード切り離し操作を拒否しません。

DR のボード追加操作が正常に終わると、追加されたボードの CPU デバイスは自動的にシステム操作に組み込まれます。

## SPARC: メモリーに対する DR クラスタリング

DR では、次の 2 種類のメモリーを考慮してください。

- カーネルメモリーケージ
- カーネル以外のメモリーケージ

これらの違いはその使用方法だけであり、実際のハードウェアは同じものです。カーネルメモリーケージとは、Solaris オペレーティングシステムが使用するメモリーのことで、Sun Cluster ソフトウェアは、カーネルメモリーケージを含むボードに対するボード切り離し操作をサポートしていないため、このような操作を拒否します。DR のボード切り離し操作がカーネルメモリーケージ以外のメモリーに関連するものである場合、Sun Cluster ソフトウェアはこの操作を拒否しません。メモリーに関連する DR のボード追加操作が正常に終わると、追加されたボードのメモリーは自動的にシステム操作に組み込まれます。

## SPARC: ディスクドライブとテープドライブに対する DR クラスタリング

Sun Cluster は、主ホストのアクティブなドライブに対する動的再構成 (Dynamic Reconfiguration、DR) のボード切り離し操作を拒否します。DR のボード切り離し操作を実行できるのは、主ホストのアクティブでないドライブと、二次ホストの任意のドライブです。DR 操作が終了すると、クラスタのデータアクセスが前と同じように続けられます。

---

注 - Sun Cluster は、定足数デバイスの使用に影響を与える DR 操作を拒否します。定足数デバイスの考慮事項と、定足数デバイスに対する DR 操作の実行手順については、108 ページの「SPARC: 定足数デバイスに対する DR クラスタリング」を参照してください。

---

これらの操作の詳細な実行手順については、『Sun Cluster のシステム管理 (Solaris OS 版)』の「定足数デバイスへの動的再構成」を参照してください。

## SPARC: 定足数デバイスに対する DR クラスタリング

DR のボード切り離し操作が、定足数デバイスとして構成されているデバイスへのインタフェースを含むボードに関連する場合、Sun Cluster ソフトウェアはこの操作を拒否します。Sun Cluster ソフトウェアはまた、この操作によって影響を受ける定足数デバイスを特定します。定足数デバイスとしてのデバイスに対して DR のボード切り離し操作を行う場合は、まずそのデバイスを無効にする必要があります。

定足数の詳細な管理手順については、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の第6章「[定足数の管理](#)」を参照してください。

## SPARC: クラスタインターコネクティングインターフェースに対する DR クラスタリング

DR のボード切り離し操作が、アクティブなクラスタインターコネクティングインターフェースを含むボードに関連する場合、Sun Cluster ソフトウェアはこの操作を拒否します。Sun Cluster ソフトウェアはまた、この操作によって影響を受けるインターフェースを特定します。DR 操作を成功させるためには、Sun Cluster 管理ツールを使用して、アクティブなインターフェースを無効にしておく必要があります。



注意 - Sun Cluster ソフトウェアでは、各クラスタノードは、ほかのすべてのクラスタノードへの有効なパスを、少なくとも1つ、持っておく必要があります。したがって、クラスタ内の個々の Solaris ホストへの最後のパスをサポートするプライベートインターコネクティングインターフェースは無効にしないでください。

これらの操作の詳細な実行方法については、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の「[クラスタインターコネクティングの管理](#)」を参照してください。

## SPARC: パブリックネットワークインターフェースに対する DR クラスタリング

DR のボード切り離し操作が、アクティブなパブリックネットワークインターフェースを含むボードに関連する場合、Sun Cluster ソフトウェアはこの操作を拒否します。Sun Cluster ソフトウェアはまた、この操作によって影響を受けるインターフェースを特定します。アクティブなネットワークインターフェースが存在するボードを切り離す前に、まず、`if_mpadm` コマンドを使って、そのインターフェース上のすべてのトラフィックを、同じマルチパスグループの正常なほかのインターフェースに切り替える必要があります。



注意 - 無効にしたネットワークアダプタに対する DR 切り離し操作中に、残りのネットワークアダプタで障害が発生すると、可用性が影響を受けます。これは、DR 操作の間は、残りのネットワークアダプタのフェイルオーバー先が存在しないためです。

パブリックネットワークインターフェースに対する DR 切り離し操作の詳細な実行手順については、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の「[パブリックネットワークの管理](#)」を参照してください。



## よくある質問

---

この章では、Sun Cluster 製品に関してもっとも頻繁に寄せられる質問に対する回答を示します。

回答は、トピックにより次のように構成されています。

- 111 ページの「高可用性に関する FAQ」
- 112 ページの「ファイルシステムに関する FAQ」
- 114 ページの「ボリューム管理に関する FAQ」
- 114 ページの「データサービスに関する FAQ」
- 115 ページの「パブリックネットワークに関する FAQ」
- 116 ページの「クラスタメンバーに関する FAQ」
- 117 ページの「クラスタ記憶装置に関する FAQ」
- 117 ページの「クラスタインターコネクトに関する FAQ」
- 118 ページの「クライアントシステムに関する FAQ」
- 118 ページの「管理コンソールに関する FAQ」
- 119 ページの「端末集配信装置とシステムサービスプロセッサに関する FAQ」

### 高可用性に関する FAQ

質問: 可用性の高いシステムとは何ですか。

回答: Sun Cluster ソフトウェアでは、高可用性 (High Availability、HA) を、クラスタがアプリケーションを実行し続けることができる能力であると定義しています。通常ならばホストシステムが使用できなくなるような障害が発生しても、高可用性アプリケーションは動作し続けます。

質問: クラスタが高可用性を提供するプロセスは何ですか。

回答: クラスタフレームワークは、フェイルオーバーとして知られるプロセスによって可用性の高い環境を提供します。フェイルオーバーとは、障害の発生したノードからクラスタ内の別の動作可能ノードにデータサービスリソースを移行するために、クラスタによって実行される一連のステップです。

質問: フェイルオーバーデータサービスとスケーラブルデータサービスの違いは何ですか。

回答: 高可用性データサービスには、次の2つの種類があります。

- フェイルオーバー
- スケーラブル

フェイルオーバーデータサービスとは、アプリケーションが一度に1つのクラスタ内の主ノードだけで実行されることを示します。ほかのノードは、ほかのアプリケーションを実行できますが、各アプリケーションは単一のノードでのみ実行されます。主ノードで障害が発生した場合、そのノードで実行中のアプリケーションは、別のノードにフェイルオーバーします。アプリケーションは実行を継続します。

スケーラブルデータサービスは、アプリケーションを複数のノードに広げて、単一の論理サービスを作成します。スケーラブルサービスは、実行されるクラスタ全体のノードとプロセッサの数を強化します。

クラスタへの物理インタフェースは、アプリケーションごとに1つのノードに設定されます。このノードを広域インタフェース(Global Interface、GIF)ノードといいます。クラスタには、複数のGIFノードが存在することがあります。個々のGIFには、スケーラブルサービスから使用する1つまたは複数の論理インタフェースがあります。この論理インタフェースを「広域インタフェース」と呼びます。GIFノードは、特定のアプリケーションに対するすべての要求を広域インタフェースを介して受け取り、それらを、そのアプリケーションサーバーが動作している複数のノードに振り分けます。GIFノードに障害が発生すると、広域インタフェースは別のノードにフェイルオーバーされます。

アプリケーションが実行されているノードに障害が発生すると、アプリケーションは別のノードで実行を続けますが、障害が発生したノードがクラスタに戻るまで多少のパフォーマンス低下が生じます。このプロセスは、障害が発生したノードがクラスタに戻るまで続けられます。

## ファイルシステムに関するFAQ

質問: クラスタ内の1つまたは複数のSolarisホストを高可用性NFSサーバーとして実行し、ほかのSolarisホストをクライアントとして実行できますか。

回答: 実行できません。ループバックマウントは行わないでください。

質問: リソースグループマネージャーの制御下でないアプリケーションにクラスタファイルシステムを使用できますか。

回答: はい。ただし、RGMの制御下ないと、そのアプリケーションが実行されているノードに障害があった場合、そのアプリケーションを手動で再起動する必要があります。



質問: クラスタファイルシステムは、必ず、/global ディレクトリの下にマウントポイントが必要ですか。

回答: いいえ。ただし、クラスタファイルシステムを /global などの同一のマウントポイントのもとに置くと、これらのファイルシステムの構成と管理が簡単になります。

質問: クラスタファイルシステムを使用した場合と NFS ファイルシステムをエクスポートした場合の違いは何ですか。

回答: 次のように、いくつかの違いがあります。

1. クラスタファイルシステムは広域デバイスをサポートします。NFS は、デバイスへの遠隔アクセスをサポートしません。
2. クラスタファイルシステムには広域名前空間があります。したがって、必要なのは1つのマウントコマンドだけです。これに対し、NFS では、ファイルシステムを各ホストにマウントする必要があります。
3. クラスタファイルシステムは、NFS よりも多くの場合でファイルをキャッシュします。たとえば、複数のノードからファイルにアクセスしている場合(たとえば、読み取り、書き込み、ファイルロック、非同期入出力などのために)、クラスタファイルシステムはファイルをキャッシュします。
4. クラスタファイルシステムは、リモート DMA とゼロコピー機能を提供する、将来の高速クラスタインターコネクトを利用するよう作られています。
5. クラスタファイルシステムのファイルの属性を (chmod などを使用して) 変更すると、変更内容はすべてのノードでただちに反映されます。エクスポートされた NFS ファイルシステムでは、この処理に時間がかかる場合があります。

質問: 私のクラスタノードには、/global/.devices/node@nodeID というファイルシステムがあります。このファイルシステムにデータを格納すると、これらのデータは高可用性および広域になりますか。

回答: 広域デバイス名前空間が格納されているこれらのファイルシステムは、一般的な使用を目的としたものではありません。これらのファイルシステムは広域的ですが、広域的にアクセスされることはありません。各ノードは、自身の広域デバイス名前空間にしかアクセスしません。あるノードが停止しても、ほかのノードがこのノードに代わってこの名前空間にアクセスすることはできません。これらのファイルシステムは、高可用性を備えてはいません。したがって、高可用性や広域属性を与えたいデータをこれらのファイルシステムに格納すべきではありません。

## ボリューム管理に関する FAQ

質問: すべてのディスクデバイスをミラー化する必要がありますか。

回答: ディスクデバイスの可用性を高くするには、それをミラー化するか、RAID-5 ハードウェアを使用する必要があります。すべてのデータサービスは、可用性の高いディスクデバイスか、可用性の高いディスクデバイスにマウントされたクラスタファイルシステムのどちらかを使用する必要があります。このような構成にすることで、単一のディスク障害に耐えることができます。

質問: ローカルディスク (起動ディスク) に対してあるボリュームマネージャーを使用し、多重ホストディスクに対して別のボリュームマネージャーを使用することはできますか。

回答: この構成をサポートするには、Solaris Volume Manager ソフトウェアでローカルディスクを管理し、Veritas Volume Manager で多重ホストディスクを管理する必要があります。これ以外の組み合わせではサポートされません。

## データサービスに関する FAQ

質問: どの Sun Cluster データサービスが利用できますか。

回答: サポートされているデータサービスのリストは、『[Sun Cluster リリースノートご使用にあたって \(Solaris OS 版\)](#)』に記載されています。

質問: Sun Cluster データサービスによってサポートされているアプリケーションのバージョンは何ですか。

回答: サポートされているアプリケーションのバージョンのリストは、『[Sun Cluster リリースノートご使用にあたって \(Solaris OS 版\)](#)』に記載されています。

質問: 独自のデータサービスを作成できますか。

回答: はい。詳細は、『[Sun Cluster データサービス開発ガイド \(Solaris OS 版\)](#)』の第 11 章「[DSDL API 関数](#)」を参照してください。

質問: ネットワークリソースを作成する場合、IP アドレスで指定するのですか。それともホスト名で指定するのですか。

回答: ネットワークリソースを指定する場合には、IP アドレスではなく、UNIX のホスト名を使用することを推奨します。

質問: ネットワークリソースを作成する場合に、論理ホスト名 (LogicalHostname リソース) または共有アドレス (SharedAddress リソース) を使用した場合の違いは何ですか。

回答: Sun Cluster HA for NFS の場合を除き、Failover モードリソースグループの LogicalHostname リソースを使用するようにマニュアルが推奨している場合

、SharedAddress リソースと LogicalHostname リソースは同様に使用できます。SharedAddress リソースを使用すると、クラスタネットワークングソフトウェアが LogicalHostname ではなく、SharedAddress に合わせて構成されているために、多少のオーバーヘッドが生じます。

SharedAddress リソースを使用する利点は、スケーラブルデータサービスとフェイルオーバーデータサービスを両方構成して、クライアントが同じホスト名で両方のサービスにアクセスするときに分かります。この場合、SharedAddress リソースは、フェイルオーバーアプリケーションリソースとともに、1つのリソースグループに格納されます。スケーラブルサービスリソースは、異なるリソースグループに格納され、SharedAddress リソースを使用するように構成されます。次に、スケーラブルサービスとフェイルオーバーサービスは両方とも、SharedAddress リソースに構成されている同じホスト名とアドレスのセットを使用します。

## パブリックネットワークに関するFAQ

質問: Sun Cluster ソフトウェアはどのパブリックネットワークアダプタをサポートしていますか。

回答: 現在、Sun Cluster ソフトウェアは、Ethernet (10/100BASE-T および 1000BASE-SX Gb) パブリックネットワークアダプタをサポートしています。今後新しいインタフェースがサポートされる可能性があるため、最新情報については、ご購入先に確認してください。

質問: フェイルオーバーでの MAC アドレスの役割は何ですか。

回答: フェイルオーバーが発生すると、新しいアドレス解決プロトコル (ARP) パケットが生成されて伝送されます。これらの ARP パケットには、新しい MAC アドレス (ホストの処理が継続される新しい物理アダプタのアドレス) と古い IP アドレスが含まれます。ネットワーク上の別のマシンがこれらのパケットの1つを受信した場合は、そのマシンは自身の ARP キャッシュから古い MAC-IP マッピングをフラッシングして、新しいマッピングを使用します。

質問: Sun Cluster ソフトウェアは local-mac-address?=true という設定をサポートしますか。

回答: はい。実際、IP ネットワークマルチパスでは local-mac-address? を true に設定する必要があります。

local-mac-address を設定するには、SPARC ベースのクラスタでは OpenBootPROM の ok プロンプトで eeprom コマンドを使用します。詳細は、[eeprom\(1M\)](#) のマニュアルページを参照してください。x86 ベースのクラスタでは、BIOS のブート後に SCSI ユーティリティを起動して設定します。

質問: IP ネットワークマルチパスがアダプタのスイッチオーバーを実行するとき、どれくらいの遅延がありますか。

回答: この遅延は数分に及ぶことがあります。これは、IP ネットワークマルチパススイッチオーバーが実行されるときに、余分な ARP ブロードキャストが送信されるためです。ただし、クライアントとクラスタ間のルーターは、必ずしもこの余分な ARP を使用するわけではありません。したがって、ルーター上のこの IP アドレスに対応する ARP キャッシュがタイムアウトするまでは、エントリが古い MAC アドレスを使用してしまう可能性があります。

質問: ネットワークアダプタの障害の検出にはどの程度の時間が必要ですか。

回答: デフォルトの障害検出時間は 10 秒です。アルゴリズムは障害をこの時間内に検出しようとはしますが、実際の時間はネットワークの負荷によって異なります。

## クラスタメンバーに関する FAQ

質問: すべてのクラスタメンバーが同じ root パスワードを持つ必要がありますか。

回答: 各クラスタメンバーに同じ root パスワードを設定する必要はありません。ただし、同じ root パスワードをすべてのノードに使用すると、クラスタの管理を簡略化できます。

質問: ノードが起動される順序は重要ですか。

回答: ほとんどの場合、重要ではありません。しかし、起動順序は `amnesia` を防ぐために重要です。たとえば、ノード 2 が定足数デバイスの所有者であり、ノード 1 が停止してノード 2 を停止させた場合は、ノード 2 を起動してからノード 1 を起動する必要があります。この順序によって、古いクラスタ構成情報を持つノードを誤って起動するのを防ぐことができます。

質問: クラスタノードのローカルディスクをミラー化する必要がありますか。

回答: はい。このミラー化は必要条件ではありませんが、クラスタノードのディスクをミラー化すると、ノードを停止させる非ミラー化ディスクの障害を防止できます。ただし、クラスタノードのローカルディスクをミラー化すると、システム管理の負荷が増えます。

質問: クラスタメンバーのバックアップの注意点は何かですか。

回答: クラスタには、いくつかのバックアップ方式を使用できます。1つの方法としては、テープドライブまたはライブラリが接続された1つのホストをバックアップノードとして設定します。さらに、クラスタファイルシステムを使用してデータをバックアップします。このホストは共有ディスクには接続しないでください。

データのバックアップと復元方法についての詳細は、『[Sun Cluster のシステム管理 \(Solaris OS 版\)](#)』の第 11 章「[クラスタのバックアップと復元](#)」を参照してください。

質問: ノードが、二次ノードとして使用できる状態にあるのはいつですか。

回答: Solaris 9 OS

再起動後にノードがログインプロンプトを表示しているときです。

Solaris 10 OS

multi-user-server マイルストーンが動作している場合、ノードは二次ノードとして使用できる状態にあります。

```
# svcs -a | grep multi-user-server:default
```

## クラスタ記憶装置に関するFAQ

質問: 多重ホスト記憶装置の可用性を高めるものは何ですか。

回答: 多重ホスト記憶装置は、ミラー化(またはハードウェアベースの RAID-5 コントローラ)によって、単一のディスクが失われても存続できるという点で高可用性です。多重ホスト記憶装置には複数のホスト接続があるため、接続先の単一の Solaris ホストが失われても耐えることができます。さらに、各ホストから、接続されている記憶装置への冗長パスは、ホストバスアダプタやケーブル、ディスクコントローラの障害に対する備えとなります。

## クラスタインターコネクトに関するFAQ

質問: Sun Cluster ソフトウェアがサポートするクラスタインターコネクトは何ですか。

回答: 現在のところ、Sun Cluster ソフトウェアは次のクラスタインターコネクトをサポートします。

- Ethernet (100BASE-T Fast Ethernet と 1000BASE-SX Gb)。SPARC ベースのクラスタと x86 ベースのクラスタの両方。
- Infiniband。SPARC ベースのクラスタと x86 ベースのクラスタの両方。
- SCI。SPARC ベースのクラスタのみ。

質問: 「ケーブル」とトランスポート「パス」の違いは何ですか。

回答: クラスタトランスポートケーブルは、トランスポートアダプタとスイッチを使用して構成されます。ケーブルは、アダプタやスイッチをコンポーネント対コンポーネントとして結合します。クラスタポロジマネージャーは、利用可能なケーブルを使用し、ホスト間にエンドツーエンドのトランスポートパスを構築します。ただし、ケーブルとトランスポートパスが 1 対 1 で対応しているわけではありません。ケーブルは、管理者によって静的に「有効」または「無効」にされます。ケーブルには、「状態」(有効または無効)はありますが、「ステータス」はありません。無効になっているケーブルは、構成されていないのと同じことです。無効なケーブル

をトランスポートパスとして使用することはできません。ケーブルは検査できないため、その状態は不明です。ケーブルの状態を取得するには、`cluster status` コマンドを使用します。

トランスポートパスは、クラスタトポロジマネージャーによって動的に確立されます。トランスポートパスの「ステータス」はトポロジマネージャーによって決められますが、パスは「オンライン」または「オフライン」のステータスを持つことができます。トランスポートパスのステータスを取得するには、`clinterconnect status` コマンドを使用します。詳細は、`clinterconnect(1CL)` のマニュアルページを参照してください。

次のような2ホストクラスタがあるとします。これには、4つのケーブルが使用されています。

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

これらの4つのケーブルを使用して設定できるトランスポートパスには、次の2つがあります。

```
node1:adapter0      to node2:adapter0
node2:adapter1      to node2:adapter1
```

## クライアントシステムに関する FAQ

質問: クラスタでの使用における特殊なクライアントの要求や制約について考慮する必要がありますか。

回答: クライアントシステムは、ほかのサーバーに接続する場合と同様にクラスタに接続します。データサービスアプリケーションによっては、クライアント側ソフトウェアをインストールするか、別の構成変更を行なって、クライアントがデータサービスアプリケーションに接続できるようにしなければならないこともあります。クライアント側の構成要件についての詳細は、『Sun Cluster データサービスの計画と管理 (Solaris OS 版)』の第1章「Sun Cluster データサービスの計画」を参照してください。

## 管理コンソールに関する FAQ

質問: Sun Cluster ソフトウェアには管理コンソールが必要ですか。

回答: はい。

質問: 管理コンソールをクラスタ専用にする必要がありますか、または別の作業に使用することができますか。

回答: Sun Cluster ソフトウェアでは専用の管理コンソールは必要ありませんが、専用の管理コンソールを使用すると、次のような利点があります。

- コンソールと管理ツールを同じマシンにまとめることで、クラスタ管理を一元化できます。
- ハードウェアサービスプロバイダによる問題解決が迅速に行われます。

質問: 管理コンソールはクラスタの近く (たとえば同じ部屋) に配置する必要がありますか。

回答: ハードウェアの保守担当者に確認してください。プロバイダによっては、コンソールをクラスタの近くに置くことを要求するところもあります。コンソールを同じ部屋に配置する必要性は、技術的にはありません。

質問: 距離の条件をすべて満たしている場合、1台の管理コンソールが複数のクラスタにサービスを提供できますか。

回答: はい。複数のクラスタを1台の管理コンソールから制御できます。また、1台の端末集配信装置 (コンセントレータ) をクラスタ間で共有することもできます。

## 端末集配信装置とシステムサービスプロセッサに関する FAQ

質問: Sun Cluster ソフトウェアは端末集配信装置を必要としますか。

回答: Sun Cluster 3.0 から、端末集配信装置は必要はありません。Sun Cluster 2.2 とは異なり、Sun Cluster 3.0、Sun Cluster 3.1、および Sun Cluster 3.2 では端末集配信装置が必要ありません。Sun Cluster 2.2 では、障害による影響防止に端末集配信装置が必要でした。

質問: ほとんどの Sun Cluster サーバーは端末集配信装置を使用していますが、Sun Enterprise E1000 サーバーが使用していないのはなぜですか。どうすればよいでしょうか。

回答: 端末集配信装置は、ほとんどのサーバーで効率的なシリアル-Ethernet コンバータです。端末集配信装置のコンソールポートはシリアルポートです。Sun Enterprise E1000 サーバーはシリアルポートを持っていません。システムサービスプロセッサ (System Service Processor、SSP) は Ethernet または jtag ポートを介したコンソールです。Sun Enterprise E1000 サーバーの場合、コンソールには常に SSP を使用します。

質問: 端末集配信装置を使用する場合の利点は何ですか。

回答: 端末集配信装置を使用すると、コンソールレベルのアクセス権が各 Solaris ホストに提供され、ネットワーク上の任意の場所にあるリモートマシンから各 Solaris ホ

ストにアクセスできます。このアクセス権は、そのホストが SPARC ベースのホスト上にある OpenBoot PROM (OBP) である場合でも、x86 ベースのホスト上にある起動サブシステムである場合でも提供されます。

質問: Sun がサポートしていない端末集配信装置を使用する場合に注意する点は何ですか。

回答: Sun がサポートする端末集配信装置とほかのコンソールデバイスの主な違いは、Sun の端末集配信装置には特殊なファームウェアがあるという点です。このファームウェアは、端末集配信装置がコンソールに対して起動時にブ레이크を送信するのを防ぎます。コンソールデバイスがブ레이크 (またはコンソールがブ레이크と解釈する可能性があるシグナル) を送信する可能性がある場合、そのブ레이크によってホストが停止されてしまうので注意してください。

質問: Sun がサポートする端末集配信装置がロックされた場合、再起動せずに、そのロックを解除できますか。

回答: はい。リセットする必要があるポート番号を書きとめて、次のコマンドを入力してください。

```
telnet tc
Enter Annex port name or number: cli
annex: su -
annex# admin
admin : reset port-number
admin : quit
annex# hangup
#
```

Sun がサポートする端末集配信装置を構成および管理する方法についての詳細は、次のマニュアルを参照してください。

- 『Sun Cluster のシステム管理 (Solaris OS 版)』の「Sun Cluster の管理の概要」
- 『Sun Cluster 3.1 - 3.2 Hardware Administration Manual for Solaris OS』の第 2 章「Installing and Configuring the Terminal Concentrator」

質問: 端末集配信装置自体に障害が発生した場合はどのようにしたらいいですか。別の装置を用意しておく必要がありますか。

回答: ありません。端末集配信装置に障害が発生しても、クラスタの可用性はまったく失われません。ただし端末集配信装置が再び機能するまでは、ホストコンソールに接続できなくなります。



質問: 端末集配信装置を使用する場合に、セキュリティーはどのように制御しますか。

回答: 通常、端末集配信装置は、ほかのクライアントアクセスに使用されるネットワークではなく、システム管理者が使用する小規模なネットワークに接続されています。この特定のネットワークに対するアクセスを制限することでセキュリティーを制御できます。

質問: SPARC: テープドライブやディスクドライブに対して動的再構成をどのように使用するのですか。

回答: 次の手順を実行します。

- ディスクドライブやテープドライブが、アクティブなデバイスグループに属しているかどうかを確認します。ドライブがアクティブなデバイスグループに属していない場合は、そのドライブに対して DR 切り離し操作を行うことができます。
- DR 切り離し操作によってアクティブなディスクドライブやテープドライブに影響がある場合には、システムは操作を拒否し、操作によって影響を受けるドライブを特定します。そのドライブがアクティブなデバイスグループに属している場合は、108 ページの「SPARC: ディスクドライブとテープドライブに対する DR クラスタリング」に進みます。
- ドライブが主ノードのコンポーネントであるか、二次ノードのコンポーネントであるかを確認します。ドライブが二次ノードのコンポーネントである場合は、そのドライブに対して DR 切り離し操作を行うことができます。
- ドライブが主ノードのコンポーネントである場合は、主ノードと二次ノードを切り替えてから、そのデバイスに対して DR 切り離し操作を行う必要があります。



注意 - 二次ノードに対して DR 操作を行っているときに現在の主ノードに障害が発生すると、クラスタの可用性が損なわれます。これは、新しい二次ノードが提供されるまでは、主ノードのフェイルオーバー先が存在しないためです。



# 索引

---

## 数字・記号

2つのホストにわたる単一クラスタポロジ, 37  
2つのホストにわたる複数のクラスタポロジ, 38-39

## A

amnesia, 63  
API, 80-82, 85

## C

CCP, 30  
CCR, 50  
CD-ROMドライブ, 27  
clprivnetドライブ, 83  
Cluster Control Panel, 30  
CMM, 48  
    フェイルファースト機構  
    「フェイルファースト」も参照  
CPUの制御, 93  
CPU時間, 95-104  
CPU, 制御, 93

## D

/dev/global/ 名前空間, 55-56  
DID, 51  
DR, 「動的再構成」を参照  
DSDL API, 85

## E

E10000, 「Sun Enterprise E10000」を参照

## F

FAQ, 111-121  
    クライアントシステム, 118  
    クラスタインターコネクト, 117-118  
    クラスタメンバー, 116-117  
    クラスタ記憶装置, 117  
    ステムサービスプロセッサ, 119-121  
    データサービス, 114-115  
    パブリックネットワーク, 115-116  
    ファイルシステム, 112-113  
    ボリューム管理, 114  
    管理コンソール, 118-119  
    高可用性, 111-112  
    端末集配信装置, 119-121

## G

/global マウントポイント, 56-59, 112-113

## H

HASStoragePlus リソースタイプ, 58-59, 84-87  
HA, 「高可用性」を参照

**I**

- ID
  - デバイス, 51
  - ノード, 55
- in.mpathd デーモン, 105
- IP アドレス, 114-115
- IP ネットワークマルチパス, 105-106
  - フェイルオーバー時間, 115-116
- IPMP, 「IP ネットワークマルチパス」を参照

**L**

- local\_mac\_address, 115-116
- LogicalHostname リソースタイプ, 「論理ホスト名」を参照

**M**

- MAC アドレス, 115-116

**N**

- N+1(星形)トポロジ, 33, 43
- N\*N(スケーラブル)トポロジ, 34
- NFS, 59
- NTP, 46-47
- numsecondaries プロパティ, 53

**O**

- Oracle Parallel Server, 「Oracle Real Application Clusters」を参照
- Oracle Real Application Clusters, 81

**P**

- per-host アドレス, 82-83
- preferenced プロパティ, 53
- pure サービス, 77

**R**

- Resource Group Manager, 「RGM」を参照
- Resource\_project\_name プロパティ, 97-98
- RG\_project\_name プロパティ, 97-98
- RGM, 74, 84-87, 95-104
- RMAPL, 85
- root パスワード, 116-117

**S**

- scha\_cluster\_get コマンド, 83
- scha\_privatelink\_hostname\_node 引数, 83
- scsi-initiator-id プロパティ, 26
- SCSI, 多重イニシエータ, 25-26
- SharedAddress リソースタイプ, 「共有アドレス」を参照
- SMF デーモン svc.startd, 92
- SMF, 「サービス管理機能 (Service Management Facility、SMF)」を参照
- Solaris Resource Manager, 95-104
  - フェイルオーバーシナリオ, 99-104
  - 仮想メモリー制限の設定, 98-99
  - 構成条件, 97-98
- Solaris Volume Manager, 多重ホストデバイス, 25
- Solaris プロジェクト, 95-104
- split brain, 63
- SSP, 「システムサービスプロセッサ」を参照
- sticky サービス, 77
- Sun Cluster Manager, 46
  - システムリソース使用量, 94
- Sun Cluster, 「クラスタ」を参照
- Sun Enterprise E10000, 119-121
  - 管理コンソール, 30
- Sun Management Center (SunMC), 46
- SUNW.Proxy\_SMF\_failover, リソースタイプ, 91
- SUNW.Proxy\_SMF\_loadbalanced, リソースタイプ, 91
- SUNW.Proxy\_SMF\_multimaster, リソースタイプ, 91
- svc.startd, デーモン, 92
- syncdir マウントオプション, 59

**U**

- UFS, 59

## V

Veritas Volume Manager, 多重ホストデバイス, 25  
VxFS, 59

## し

## しきい値

システムリソース, 93  
テレメトリ属性, 93

## よ

よくある質問, 「FAQ」を参照

## ア

アダプタ, 「ネットワーク、アダプタ」を参照  
アプリケーション, 「データサービス」を参照  
アプリケーション開発, 45-109  
アプリケーション通信, 82-83  
アプリケーション配布, 67

## イ

## インタフェース

「ネットワーク、インタフェース」を参照  
管理, 46

## エ

エージェント, 「データサービス」を参照

## オ

オブジェクトタイプ, システムリソース, 92

## カ

カーネル, メモリー, 108

## ク

クライアントサーバー構成, 72

クライアントシステム

FAQ, 118

制限, 118

クラスタ

アプリケーション開発, 45-109

アプリケーション開発者, 18-19

インターコネクト

FAQ, 117-118

アダプタ, 27

インタフェース, 27

ケーブル, 28

サポートされる, 117-118

データサービス, 82-83

接続点, 28

動的再構成, 109

サービス, 16

システム管理者, 17-18

ソフトウェアコンポーネント, 23-24

データサービス, 71-79

トポロジ, 31-41, 42-43

ノード, 22-23

ハードウェア, 16, 21-30

バックアップ, 116-117

パスワード, 116-117

パブリックネットワーク, 28

パブリックネットワークインタフェース, 72

ファイルシステム

FAQ

「ファイルシステム」も参照

HASStoragePlus リソースタイプ, 58-59

使用法, 57-58

ボード切り離し, 108

メディア, 27

メンバー

FAQ, 116-117

再構成, 48

管理, 45-109

記憶装置に関する FAQ, 117

## クラスタ (続き)

- 起動順序, 116-117
- 構成, 50, 95-104
- 作業リスト, 19-20
- 時間, 46-47
- 説明, 13-15
- 目的, 13-15
- 利点, 13-15
- クラスタベアトポロジ, 31-32, 42
- クラスタメンバーシップモニター, 48
- クラスタ化サーバーモデル, 72
- クラスタ構成レポジトリ, 50

## グ

- グループ, デバイス, 51-55

## ケ

- ケーブル, トランスポート, 117-118

## コ

- コンソール
  - アクセス, 29-30
  - システムサービスプロセッサ, 29-30
  - 管理
    - FAQ, 118-119

## サ

- サーバーモデル, 72
- サービス管理機能 (Service Management Facility, SMF), 91-92

## シ

- システムサービスプロセッサ, 29-30, 30
- システムリソース
  - しきい値, 93

## システムリソース (続き)

- オブジェクトタイプ, 92
- 監視, 93
- 使用率, 92
- システムリソースの使用率, 92
- システムリソース監視, 93

## ス

- スケーラブルデータサービス, 75-76
- ステムサービスプロセッサ, FAQ, 119-121

## ソ

- ソフトウェアコンポーネント, 23-24

## ゾ

- ゾーン, 88

## テ

- テープドライブ, 27
- テレメトリ属性, システムリソース, 93

## デ

- データ, 格納, 112-113
- データサービス
  - API, 80-82
  - FAQ, 114-115
  - クラスタインターコネクト, 82-83
  - サポートされている, 114-115
  - スケーラブル, 75-76
  - フェイルオーバー, 75
  - メソッド, 74
  - ライブラリ API, 82
  - リソース, 84-87
  - リソースグループ, 84-87
  - リソースタイプ, 84-87

## データサービス (続き)

- 開発, 80-82
- 構成, 95-104
- 高可用, 48
- 障害モニター, 79

## デーモン, svc.startd, 92

## ディスク

- SCSI デバイス, 25-26
- ローカル, 26-27, 50-51, 55-56
  - ボリューム管理, 114
  - ミラー化, 116-117
- 広域デバイス, 50-51, 55-56
- 多重ホスト, 50-51, 51-55, 55-56
- 動的再構成, 108

## ディスクバス監視, 60-63

## デバイス

- ID, 51
- 広域, 50-51
- 多重ホスト, 25
- 定足数, 63-70

## デバイスグループ, 51-55

- フェイルオーバー, 52
- プロパティの変更, 53-55
- 主所有権, 53-55
- 多重ポート, 53-55

## ト

## トポロジ

- Logical Domains: 2つのホストにわたる単一クラスタ, 37
- Logical Domains: 2つのホストにわたる複数のクラスタ, 38-39
- Logical Domains: ボックス内クラスタトポロジ, 35
- Logical Domains: 冗長 I/O ドメイン, 40-41
- N+1 (星形), 33, 43
- N\*N (スケラブル), 34
- クラスタペア, 31-32, 42
- ペア +N, 32-33

## ド

ドライバ, デバイス ID, 51

## ネ

## ネットワーク

- アダプタ, 28, 105-106
- インタフェース, 28, 105-106
- パブリック
  - FAQ, 115-116
  - IP ネットワークマルチパス, 105-106
  - インタフェース, 115-116
  - 動的再構成, 109
- プライベート, 23
- リソース, 72, 84-87
- 共有アドレス, 72
- 負荷均衡, 77-79
- 論理ホスト名, 72

## ノ

## ノード

- nodeID, 55
- バックアップ, 116-117
- 起動順序, 116-117
- 広域インタフェース, 73
- 主, 53-55, 73
- 二次, 53-55, 73

## ハ

## ハードウェア

- 「ディスク」も参照
- 「記憶装置」も参照
- クラスタインターコネクトコンポーネント, 27
- 動的再構成, 106-109

## バ

バックアップノード, 116-117

## パ

パス, トランスポート, 117-118  
パスワード, root, 116-117  
パニック, 48-50  
パブリックネットワーク, 「ネットワーク、パブリック」を参照  
パラレルデータベース構成, 23

## フ

ファイルシステム  
FAQ, 112-113  
NFS, 59, 112-113  
syncdir マウントオプション, 59  
UFS, 59  
VxFS, 59  
クラスタ, 56-59, 112-113  
データ記憶装置, 112-113  
マウント, 56-59, 112-113  
ローカル, 58-59  
広域  
「ファイルシステム、クラスタ」を参照  
高可用性, 112-113  
使用法, 57-58  
ファイルロック, 56  
フェイルオーバー  
シナリオ, Solaris Resource Manager, 99-104  
データサービス, 75  
デバイスグループ, 52  
フェイルバック, 79  
フェイルファースト, 48-50  
フェンシング, 49  
フレームワーク, 高可用性, 47-50

## プ

プライベートネットワーク, 23  
プロキシリソースタイプ, 91  
プロジェクト, 95-104  
プロパティ  
Resource\_project\_name, 97-98  
RG\_project\_name, 97-98  
リソース, 87

## プロパティ (続き)

リソースグループ, 87  
変更, 53-55

## ペ

ペア +N トポロジ, 32-33

## ホ

ホスト名, 72

## ボ

ボード切り離し, 動的再構成, 108  
ボックス内クラスタトポロジ, 35  
ボリューム管理  
FAQ, 114  
RAID-5, 114  
Solaris Volume Manager, 114  
Veritas Volume Manager, 114  
ローカルディスク, 114  
多重ホストディスク, 114  
多重ホストデバイス, 25  
名前空間, 55

## マ

マウント  
/global, 112-113  
syncdir, 59  
ファイルシステム, 56-59  
広域デバイス, 56-59  
マッピング, 名前空間, 56  
マルチパス, 105-106

## ミ

ミッションクリティカルなアプリケーション, 69



## メ

メディア, リムーバブル, 27  
 メモリー, 108  
 メンバーシップ, 「クラスタ、メンバー」を参照

## リ

リソース, 84-87  
   プロパティ, 87  
   状態, 85-87  
   設定値, 85-87  
 リソースグループ, 84-87  
   スケーラブル, 75-76  
   フェイルオーバー, 75  
   プロパティ, 87  
   状態, 85-87  
   設定値, 85-87  
 リソースタイプ  
   SUNW.Proxy\_SMF\_failover, 91  
   SUNW.Proxy\_SMF\_loadbalanced, 91  
   SUNW.Proxy\_SMF\_multimaster, 91  
   プロキシ, 91  
 リソース管理, 95-104  
 リムーバブルメディア, 27

## ロ

ローカルディスク, 26-27  
 ローカルファイルシステム, 58-59  
 ローカル名前空間, 56

## 回

回復  
   フェイルバック設定, 79  
   障害検出, 47

## 開

開発者, クラスタアプリケーション, 18-19

## 監

監視  
   オブジェクトタイプ, 92  
   システムリソース, 93  
   テレメトリ属性, 93  
   ディスクバス, 60-63

## 管

管理, クラスタ, 45-109  
 管理インタフェース, 46  
 管理コンソール, FAQ, 118-119

## 記

記憶装置  
   FAQ, 117  
   SCSI, 25-26  
   動的再構成, 108

## 起

起動ディスク, 「ディスク、ローカル」を参照  
 起動順序, 116-117

## 共

共有アドレス, 72  
   スケーラブルデータサービス, 75-76  
   広域インタフェースノード, 73  
   対論理ホスト名, 114-115

## 広

広域  
   インタフェース, 73  
   スケーラブルサービス, 76  
 デバイス, 50-51, 51-55  
   マウント, 56-59  
   ローカルディスク, 26

広域 (続き)

名前空間, 50, 55-56

ローカルディスク, 26

広域インタフェースノード, 73

広域ファイルシステム, 「クラスタ、ファイルシステム」を参照

構

構成

クライアントサーバー, 72

データサービス, 95-104

パラレルデータベース, 23

レポジトリ, 50

仮想メモリーの限度, 98-99

定足数, 65-66

高

高可用, データサービス, 48

高可用性

FAQ, 111-112

フレームワーク, 47-50

時

時間, ノード間の, 46-47

時間情報プロトコル, 46-47

主

主ノード, 73

主所有権, デバイスグループ, 53-55

障

障害

フェイルバック, 79

回復, 47

検出, 47

障害モニター, 79

冗

冗長 I/O ドメイントポロジ, 40-41

属

属性, 「プロパティ」を参照

多

多重イニシエータ SCSI, 25-26

多重ホストデバイス, 25

多重ポートデバイスグループ, 53-55

単

単一サーバーモデル, 72

端

端末集配信装置, FAQ, 119-121

停

停止, 48-50

定

定足数, 63-70

デバイス, 63-70

デバイス、動的再構成, 108-109

ベストプラクティス, 66-67

構成, 65

推奨される構成, 67-68

投票数, 64-65

変則的な構成, 69

## 定足数 (続き)

望ましくない構成, 70  
要件, 65-66

## 投

投票数, 定足数, 64-65

## 動

## 動的再構成

CPU デバイス, 107  
クラスタインターコネクト, 109  
テープドライブ, 108  
ディスク, 108  
パブリックネットワーク, 109  
メモリー, 108  
説明, 107  
定足数デバイス, 108-109

## 同

同時アクセス, 23

## 二

二次ノード, 73

## 負

負荷均衡, 77-79

## 名

名前空間, 55-56, 56

## 論

論理ホスト名, 72  
フェイルオーバーデータサービス, 75  
対共有アドレス, 114-115

