

CLUSTERPRO[®] X 1.0 *for Linux*

スタートアップガイド

2007.03.30
第3版



改版履歴

版数	改版日付	内 容
1	2006/09/08	新規作成
2	2006/12/12	CLUSTERPROロゴを変更しました。 第 3 章 CLUSTERPROの動作環境 ソフトウェア「動作可能なディストリビューションとkernel」に新規対応kernel と備考欄を追加しました。 第 5 章 注意制限事項 CLUSTERPROの情報作成時「ミラーディスクのファイルシステムの選択 について」に新規対応ファイルシステムを追加しました。
3	2007/03/30	第 3 章 CLUSTERPROの動作環境 ソフトウェア「動作可能なディストリビューションとkernel」に新規対応kernel を追加しました。

免責事項

本書の内容は、予告なしに変更されることがあります。

日本電気株式会社は、本書の技術的もしくは編集上の間違い、欠落について、一切責任をおいせん。

また、お客様が期待される効果を得るために、本書に従った導入、使用および使用効果につきましては、お客様の責任とさせていただきます。

本書に記載されている内容の著作権は、日本電気株式会社に帰属します。本書の内容の一部または全部を日本電気株式会社の許諾なしに複製、改変、および翻訳することは禁止されています。

商標情報

CLUSTERPRO[®] X は日本電気株式会社の登録商標です。

FastSync[™]は日本電気株式会社の商標です。

Linuxは、Linus Torvalds氏の米国およびその他の国における、登録商標または商標です。

RPMの名称は、Red Hat, Inc.の商標です。

Intel、Pentium、Xeonは、Intel Corporationの登録商標または商標です。

Microsoft、Windowsは、米国Microsoft Corporationの米国およびその他の国における登録商標です。

Turbolinuxおよびターボリナックスは、ターボリナックス株式会社の登録商標です。

VERITAS、VERITAS ロゴ、およびその他のすべてのVERITAS 製品名およびスローガンは、VERITAS Software Corporation の商標または登録商標です。

本書に記載されたその他の製品名および標語は、各社の商標または登録商標です。

目次

はじめに	vii
対象読者と目的	vii
本書の構成	vii
CLUSTERPRO マニュアル体系	viii
本書の表記規則	ix
最新情報の入手先	x
セクション I CLUSTERPROの概要	1
第 1 章 クラスタシステムとは?	3
クラスタシステムの概要	4
HA (High Availability) クラスタ	4
共有ディスク型	5
データミラー型	7
障害検出のメカニズム	8
共有ディスク型の諸問題	8
ネットワークパーティション症状(Split-brain-syndrome)	9
クラスタリソースの引き継ぎ	9
データの引き継ぎ	9
アプリケーションの引き継ぎ	10
フェイルオーバー総括	11
Single Point of Failureの排除	12
共有ディスク	12
共有ディスクへのアクセスパス	13
LAN	14
可用性を支える運用	14
障害の監視	15
第 2 章 CLUSTERPRO の使用方法	17
CLUSTERPRO とは?	18
CLUSTERPRO の製品構成	18
CLUSTERPRO のソフトウェア構成	18
CLUSTERPRO の障害監視のしくみ	19
サーバ監視とは	19
業務監視とは	20
内部監視とは	20
監視できる障害と監視できない障害	21
サーバ監視で検出できる障害とできない障害	21
業務監視で検出できる障害とできない障害	21
フェイルオーバーのしくみ	22
フェイルオーバーリソース	23
フェイルオーバー型クラスタのシステム構成	23
共有ディスク型のハードウェア構成	26
ミラーディスク型のハードウェア構成	27
クラスタオブジェクトとは?	28
リソースとは?	29
ハートビートリソース	29
グループリソース	29
モニタリソース	30
CLUSTERPRO を始めよう!	31

最新情報の確認	31
クラスタシステムの設計	31
クラスタシステムの構築	31
クラスタシステムの運用開始後の障害対応	31
セクション II リリースノート (CLUSTERPRO 最新情報).....	33
第 3 章 CLUSTERPRO の動作環境.....	35
ハードウェア	36
スペック 36	
動作確認済ディスクインターフェイス	36
動作確認済ネットワークインターフェイス	37
ソフトウェア	38
CLUSTERPRO Serverの動作環境	38
動作可能なディストリビューションとkernel	38
必要メモリ容量とディスクサイズ	41
Builderの動作環境	42
動作確認済OS、ブラウザ	42
Java実行環境	42
必要メモリ容量/ディスク容量	42
対応するCLUSTERPROのバージョン	42
WebManagerの動作環境	43
動作確認済OS、ブラウザ	43
Java実行環境	43
必要メモリ容量/ディスク容量	43
第 4 章 最新バージョン情報	45
第 5 章 注意制限事項.....	47
システム構成検討時の注意事項	48
Builder、WebManagerの動作OSについて	48
ミラーディスクの要件について	48
共有ディスクの要件について	49
NIC Link Up/Downモニタリソース	50
ミラーリソースのwrite性能について	50
OSインストール前、OSインストール時	51
/opt/nec/clusterproのファイルシステムについて	51
ミラー用のディスクについて	51
依存するライブラリ	53
依存するドライバ	53
ミラードライバ	53
カーネルモードLANハートビートドライバ、キープアライブドライバ	53
RAWモニタリソース用のパーティション確保	53
OSインストール後、CLUSTERPROインストール前	54
通信ポート番号	54
時刻同期の設定	55
NICデバイス名について	55
共有ディスクについて	55
ミラー用のディスクについて	57
OS起動時間の調整	58
ネットワークの確認	59
ユーザ空間モニタリソース(監視方法ipmi)について	60
ユーザ空間モニタリソース(監視方法softdog)について	60
CLUSTERPROの情報作成時	61
グループリソースの非活性異常時の最終アクション	61
VxVMが使用するRAWデバイスの確認	61

ミラーディスクのファイルシステムの選択について	62
RAWモニタリソースについて	62
遅延警告割合	62
ディスクモニタリソースの監視方法TURについて	62
WebManagerの画面更新間隔について	63
LANハートビートの設定について	63
カーネルモードLANハートビートの設定について.....	63
COMハートビートの設定について.....	63
CLUSTERPRO運用後.....	64
hotplugサービスについて	64
X-Window上のファイル操作ユーティリティについて.....	64
ドライバロード時のメッセージについて	64
ipmiのメッセージについて	65
回復動作中の操作制限.....	65
コマンド編に記載されていない実行形式ファイルやスクリプトファイルについて	65
kernelページアロケートエラーのメッセージについて.....	65
ログ収集時のメッセージ	66
クラスタシャットダウン・クラスタシャットダウンリブート	66
特定サーバのシャットダウン、リブート	66
WebManagerについて	66
Builder について	67
第 6 章 アップデート手順	69
CLUSTERPRO Ver3.xからのアップデート手順	70
クラスタ構成情報のバックアップ	70
クラスタ情報の変換.....	70
3.xのアンインストール	71
X 1.0のインストール.....	71
付録	73
付録 A 用語集	75
付録 B 索引.....	79

はじめに

対象読者と目的

『CLUSTERPRO®スタートアップガイド』は、CLUSTERPRO をはじめてご使用になるユーザの皆様を対象に、CLUSTERPRO の製品概要、クラスタシステム導入のロードマップ、他マニュアルの使用方法についてのガイドラインを記載します。また、最新の動作環境情報や制限事項などについても紹介します。

本書の構成

セクション I CLUSTERPRO の概要

- 第 1 章 「クラスタシステムとは?」: クラスタシステムおよび CLUSTERPRO の概要について説明します。
- 第 2 章 「CLUSTERPRO の使用方法」: クラスタシステムの使用方法および関連情報について説明します。

セクション II リリース ノート

- 第 3 章 「CLUSTERPRO の動作環境」: 導入前に確認が必要な最新情報について説明します。
- 第 4 章 「最新バージョン情報」: CLUSTERPRO の最新バージョンについての情報を示します。
- 第 5 章 「注意制限事項」: 既知の問題と制限事項について説明します。
- 第 6 章 「アップデート手順」: 既存バージョンから最新版へのアップデート情報について説明します。

付録

- 付録 A 「用語集」
- 付録 B 「索引」

CLUSTERPRO マニュアル体系

CLUSTERPRO のマニュアルは、以下の 4 つに分類されます。各ガイドのタイトルと役割を以下に示します。

『CLUSTERPRO X スタートアップガイド』(Getting Started Guide)

すべてのユーザを対象読者とし、製品概要、動作環境、アップデート情報、既知の問題などについて記載します。

『CLUSTERPRO X インストール & 設定ガイド』(Install and Configuration Guide)

CLUSTERPRO を使用したクラスタ システムの導入を行うシステム エンジニアと、クラスタ システム導入後の保守・運用を行うシステム管理者を対象読者とし、CLUSTERPRO を使用したクラスタ システム導入から運用開始前までに必須の事項について説明します。実際にクラスタ システムを導入する際の順番に則して、CLUSTERPRO を使用したクラスタ システムの設計方法、CLUSTERPRO のインストールと設定手順、設定後の確認、運用開始前の評価方法について説明します。

『CLUSTERPRO X リファレンス ガイド』(Reference Guide)

管理者を対象とし、CLUSTERPRO の運用手順、各モジュールの機能説明、メンテナンス関連情報およびトラブルシューティング情報等を記載します。『インストール & 設定ガイド』を補完する役割を持ちます。

『CLUSTERPRO X (製品別) 管理者ガイド』(Add-on Products Administrator's Guide)

管理者を対象とし、CLUSTERPRO で用意されている関連製品について、製品概要、設定方法などの詳細情報を記載します。以下の 5 冊があります。

『Alert Service 管理者ガイド』

『Application Server Agent 管理者ガイド』

『Database Agent 管理者ガイド』

『File Server Agent 管理者ガイド』

『Internet Server Agent 管理者ガイド』

本書の表記規則

本書では、注意すべき事項、重要な事項および関連情報を以下のように表記します。

注：は、重要ではあるがデータ損失やシステムおよび機器の損傷には関連しない情報を表します。

重要：は、データ損失やシステムおよび機器の損傷を回避するために必要な情報を表します。

関連情報：は、参照先の情報の場所を表します。

また、本書では以下の表記法を使用します。

表記	使用方法	例
[] 角かっこ	コマンド名の前後 画面に表示される語 (ダイアログ ボックス、メニューなど) の前後	[スタート] をクリックします。 [プロパティ] ダイアログ ボックス
コマンドライン中の [] 角かっこ	かっこ内の値の指定が省略可能であることを示します。	clpstat -s[-h host_name]
#	Linux ユーザが、root でログインしていることを示すプロンプト	# clpcl -s -a
モノスペース フォント (courier)	パス名、コマンド ライン、システムからの出力 (メッセージ、プロンプトなど)、ディレクトリ、ファイル名、関数、パラメータ	/Linux/1.0/jpn/server/
モノスペース フォント太字 (courier)	ユーザが実際にコマンドラインから入力する値を示します。	以下を入力します。 # clpcl -s -a
モノスペース フォント (courier) 斜体	ユーザが有効な値に置き換えて入力する項目	rpm -i clusterprobuilder-<バージョン番号>-<リリース番号>.i686.rpm

最新情報の入手先

最新の製品情報については、以下のWebサイトを参照してください。

<http://www.ace.comp.nec.co.jp/CLUSTERPRO/index.html>

セクション I CLUSTERPRO の概要

このセクションでは、CLUSTERPRO の製品概要と動作環境について説明します。

- 第 1 章 クラスタシステムとは？
- 第 2 章 CLUSTERPRO の使用方法

第 1 章 クラスタシステムとは？

本章では、クラスタシステムの概要について説明します。

本章で説明する項目は以下のとおりです。

• クラスタシステムの概要	4
• HA (High Availability) クラスタ.....	4
• 障害検出のメカニズム	8
• クラスタリソースの引き継ぎ	9
• Single Point of Failureの排除	12
• 可用性を支える運用	14

クラスタシステムの概要

現在のコンピュータ社会では、サービスを停止させることなく提供し続けることが成功への重要なカギとなります。例えば、1 台のマシンが故障や過負荷によりダウンしただけで、顧客へのサービスが全面的にストップしてしまうことがあります。そうすると、莫大な損害を引き起こすだけでなく、顧客からの信用を失いかねません。

このような事態に備えるのがクラスタシステムです。クラスタシステムを導入することにより、万一のときのシステム稼働停止時間(ダウンタイム)を最小限に食い止めたり、負荷を分散させたりすることでシステムダウンを回避することが可能になります。

クラスタとは、「群れ」「房」を意味し、その名の通り、クラスタシステムとは「複数のコンピュータを一群(または複数群)にまとめて、信頼性や処理性能の向上を狙うシステム」です。クラスタシステムには様々な種類があり、以下の 3 つに分類できます。この中で、CLUSTERPRO はハイアベイラビリティクラスタに分類されます。

◆ HA (ハイ アベイラビリティ) クラスタ

通常時は一方が現用系として業務を提供し、現用系障害発生時に待機系に業務を引き継ぐような形態のクラスタです。高可用性を目的としたクラスタで、データの引継ぎも可能です。共有ディスク型、データミラー型、遠隔クラスタがあります。

◆ 負荷分散クラスタ

クライアントからの要求を適切な負荷分散ルールに従って負荷分散ホストに要求を割り当てるクラスタです。高スケーラビリティを目的としたクラスタで、一般的にデータの引継ぎはできません。ロードバランスクラスタ、並列データベースクラスタがあります。

◆ HPC(High Performance Computing)クラスタ

全てのノードの CPU を利用し、単一の業務を実行するためのクラスタです。高性能化を目的としたおり、あまり汎用性はありません。
なお、HPC の 1 つであり、より広域な範囲のノードや計算機クラスタまでを束ねた、グリッドコンピューティングという技術も近年話題に上ることが多くなっています。

HA (High Availability) クラスタ

一般的にシステムの可用性を向上させるには、そのシステムを構成する部品を冗長化し、Single Point of Failure をなくすことが重要であると考えられます。Single Point of Failure とは、コンピュータの構成要素 (ハードウェアの部品) が 1 つしかないために、その個所で障害が起きると業務が止まってしまう弱点のことを指します。HA クラスタとは、サーバを複数台使用して冗長化することにより、システムの停止時間を最小限に抑え、業務の可用性 (availability) を向上させるクラスタシステムをいいます。

システムの停止が許されない基幹業務システムはもちろん、ダウンタイムがビジネスに大きな影響を与えてしまうそのほかのシステムにおいても、HA クラスタの導入が求められています。

HA クラスタは、共有ディスク型とデータミラー型に分けることができます。以下にそれぞれのタイプについて説明します。

共有ディスク型

クラスタシステムでは、サーバ間でデータを引き継がなければなりません。このデータを共有ディスク上に置き、ディスクを複数のサーバで利用する形態を共有ディスク型といいます。

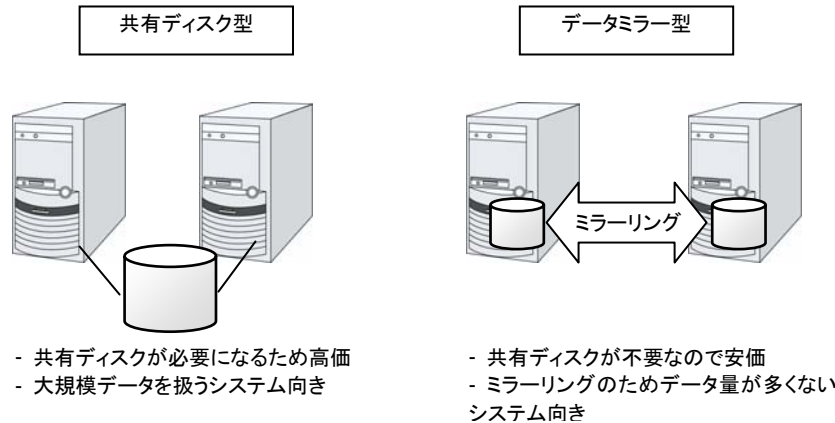


図 1-1 HAクラスタ構成図

業務アプリケーションを動かしているサーバ(現用系サーバ)で障害が発生した場合、クラスタシステムが障害を検出し、待機系サーバで業務アプリケーションを自動起動させ、業務を引き継がせます。これをフェイルオーバーといいます。クラスタシステムによって引き継がれる業務は、ディスク、IP アドレス、アプリケーションなどのリソースと呼ばれるもので構成されています。

クラスタ化されていないシステムでは、アプリケーションをほかのサーバで再起動させると、クライアントは異なる IP アドレスに再接続しなければなりません。しかし、多くのクラスタシステムでは、業務単位に仮想 IP アドレスを割り当てています。このため、クライアントは業務を行っているサーバが現用系か待機系かを意識する必要はなく、まるで同じサーバに接続しているように業務を継続できます。

データを引き継ぐためには、ファイルシステムの整合性をチェックしなければなりません。通常は、ファイルシステムの整合性をチェックするためにチェックコマンド（例えば、Linux の場合は `fsck` や `chkdsk`）を実行しますが、ファイルシステムが大きくなるほどチェックにかかる時間が長くなり、その間業務が止まってしまいます。この問題を解決するために、ジャーナリングファイルシステムなどでフェイルオーバー時間を短縮します。

業務アプリケーションは、引き継いだデータの論理チェックをする必要があります。例えば、データベースならばロールバックやロールフォワードの処理が必要になります。これらによって、クライアントは未コミットの SQL 文を再実行するだけで、業務を継続することができます。

障害からの復帰は、障害が検出されたサーバを物理的に切り離して修理後、クラスタシステムに接続すれば待機系として復帰できます。業務の継続性を重視する実際の運用の場合は、ここまでの復帰で十分な状態です。

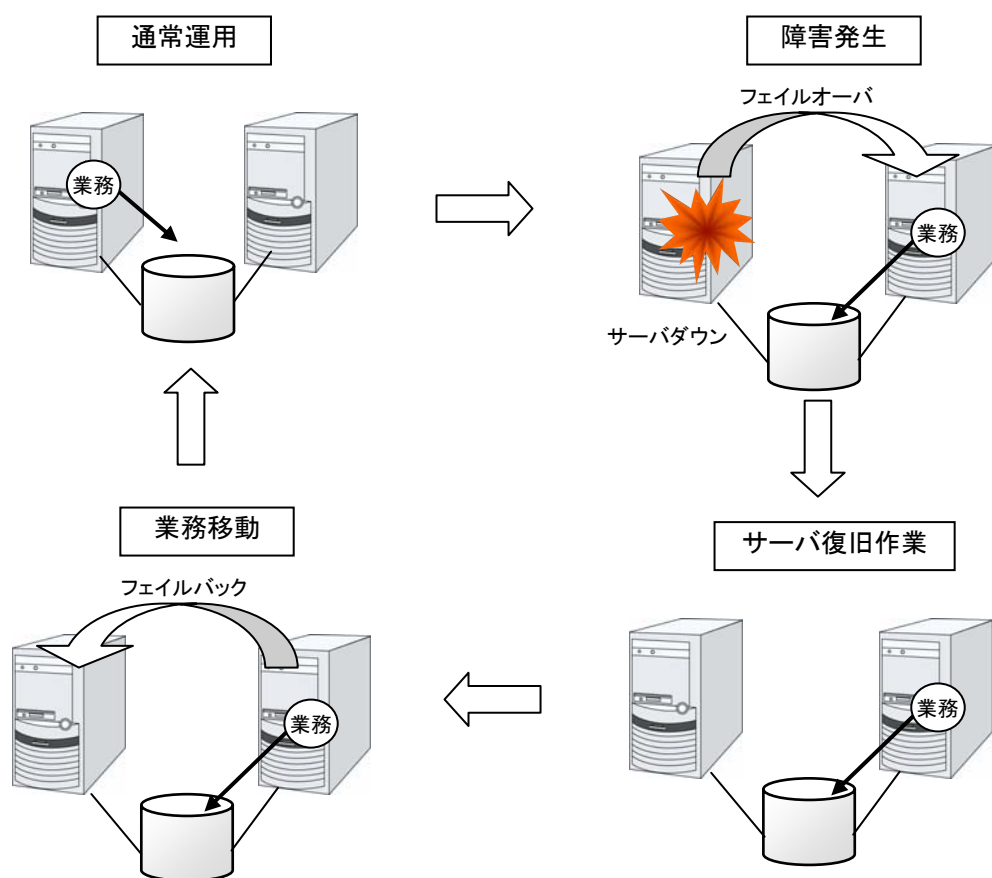


図 1-2 障害発生から復旧までの流れ

フェイルオーバー先のサーバのスペックが十分でなかったり、双方向スタンバイで過負荷になるなどの理由で元のサーバで業務を行うのが望ましい場合には、元のサーバで業務を再開するためにフェイルバックを行います。

図 1-3 のように、業務が 1 つであり、待機系では業務が動作しないスタンバイ形態を片方向スタンバイといいます。業務が 2 つ以上で、それぞれのサーバが現用系かつ待機系である形態を双方向スタンバイといいます。

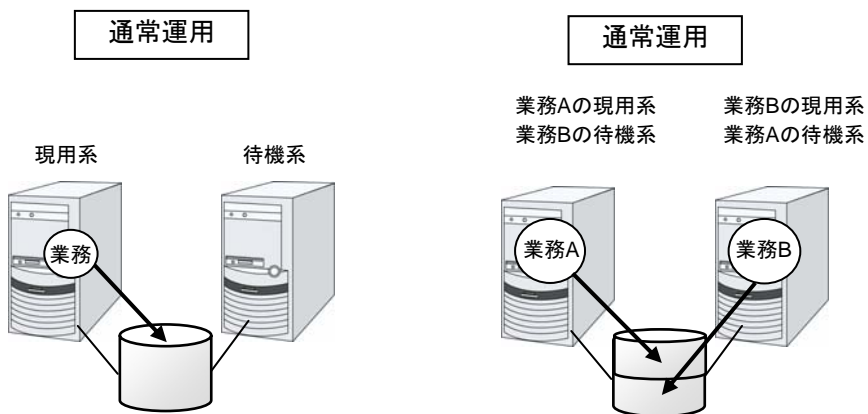


図 1-3 HA クラスターの運用形態

データミラー型

前述の共有ディスク型は大規模なシステムに適していますが、共有ディスクはおおむね高価なためシステム構築のコストが膨らんでしまいます。そこで共有ディスクを使用せず、各サーバのディスクをサーバ間でミラーリングすることにより、同じ機能をより低価格で実現したクラスタシステムをデータミラー型といいます。

しかし、サーバ間でデータをミラーリングする必要があるため、大量のデータを必要とする大規模システムには向きません。

アプリケーションからの Write 要求が発生すると、データミラーエンジンはローカルディスクにデータを書き込むと同時に、インタコネクトを通して待機系サーバにも Write 要求を振り分けます。インタコネクトとは、サーバ間をつなぐネットワークのことで、クラスタシステムではサーバの死活監視のために必要になります。データミラータイプでは死活監視に加えてデータの転送に使用することがあります。待機系のデータミラーエンジンは、受け取ったデータを待機系のローカルディスクに書き込むことで、現用系と待機系間のデータを同期します。

アプリケーションからの Read 要求に対しては、単に現用系のディスクから読み出すだけです。

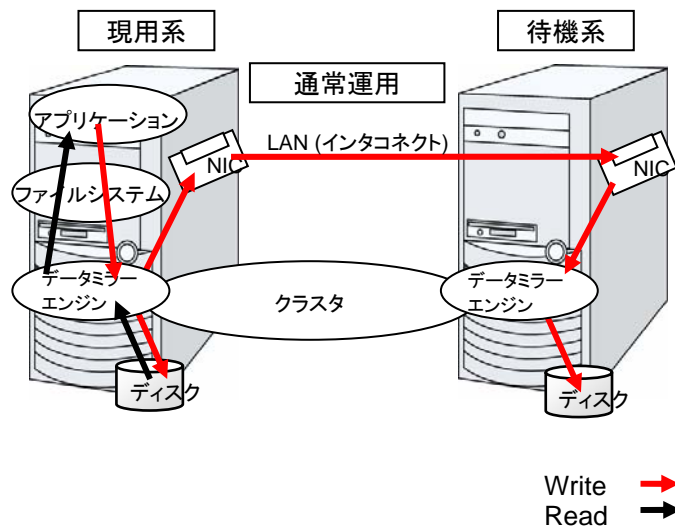


図 1-4 データミラーの仕組み

データミラーの応用例として、スナップショットバックアップの利用があります。データミラータイプのクラスタシステムは2カ所に共有のデータを持っているため、待機系のサーバをクラスタから切り離すだけで、バックアップ時間をかけることなくスナップショットバックアップとしてディスクを保存する運用が可能です。

フェイルオーバーの仕組みと問題点

ここまで、一口にクラスタシステムといってもフェイルオーバークラスタ、負荷分散クラスタ、HPC(High Performance Computing)クラスタなど、さまざまなクラスタシステムがあることを説明しました。そして、フェイルオーバークラスタは HA(High Availability)クラスタと呼ばれ、サーバそのものを多重化することで、障害発生時に実行していた業務をほかのサーバで引き継ぐことにより、業務の可用性(Availability)を向上することを目的としたクラスタシステムであることを見てきました。次に、クラスタの実装と問題点について説明します。

障害検出のメカニズム

クラスタソフトウェアは、業務継続に問題をきたす障害を検出すると業務の引き継ぎ(フェイルオーバー)を実行します。フェイルオーバー処理の具体的な内容に入る前に、簡単にクラスタソフトウェアがどのように障害を検出するか見ておきましょう。

ハートビートとサーバの障害検出

クラスタシステムにおいて、検出すべき最も基本的な障害はクラスタを構成するサーバ全てが停止してしまうものです。サーバの障害には、電源異常やメモリエラーなどのハードウェア障害や OS のパニックなどが含まれます。このような障害を検出するために、サーバの死活監視としてハートビートが使用されます。

ハートビートは、ping の応答を確認するような死活監視だけでもよいのですが、クラスタソフトウェアによっては、自サーバの状態情報などを相乗りさせて送るものもあります。クラスタソフトウェアはハートビートの送受信を行い、ハートビートの応答がない場合はそのサーバの障害とみなしてフェイルオーバー処理を開始します。ただし、サーバの高負荷などによりハートビートの送受信が遅延することも考慮し、サーバ障害と判断するまである程度の猶予時間が必要です。このため、実際に障害が発生した時間とクラスタソフトウェアが障害を検知する時間とにはタイムラグが生じます。

リソースの障害検出

業務の停止要因はクラスタを構成するサーバ全ての停止だけではありません。例えば、業務アプリケーションが使用するディスク装置や NIC の障害、もしくは業務アプリケーションそのものの障害などによっても業務は停止してしまいます。可用性を向上するためには、このようなリソースの障害も検出してフェイルオーバーを実行しなければなりません。

リソース異常を検出する手法として、監視対象リソースが物理的なデバイスの場合は、実際にアクセスしてみるという方法が取られます。アプリケーションの監視では、アプリケーションプロセスそのものの死活監視のほか、業務に影響のない範囲でサービスポートを試してみるような手段も考えられます。

共有ディスク型の諸問題

共有ディスク型のフェイルオーバークラスタでは、複数のサーバでディスク装置を物理的に共有します。一般的に、ファイルシステムはサーバ内にデータのキャッシュを保持することで、ディスク装置の物理的な I/O 性能の限界を超えるファイル I/O 性能を引き出しています。

あるファイルシステムを複数のサーバから同時にマウントしてアクセスするとどうなるでしょうか？

通常のファイルシステムは、自分以外のサーバがディスク上のデータを更新するとは考えていないので、キャッシュとディスク上のデータとに矛盾を抱えることとなり、最終的にはデータを破壊します。フェイルオーバークラスタシステムでは、次のネットワークパーティション症状などによる複数サーバからのファイルシステムの同時マウントを防ぐために、ディスク装置の排他制御を行っています。

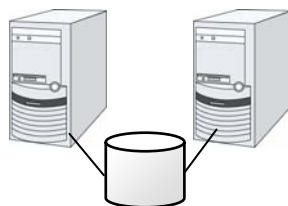


図 1-5 共有ディスクタイプのクラスタ構成

ネットワークパーティション症状(Split-brain-syndrome)

サーバ間をつなぐすべてのインタコネクトが切断されると、ハートビートによる死活監視で互いに相手サーバのダウンを検出し、フェイルオーバー処理を実行してしまいます。結果として、複数のサーバでファイルシステムを同時にマウントしてしまい、データ破壊を引き起こします。フェイルオーバークラスタシステムでは異常が発生したときに適切に動作しなければならないことが理解できると思います。

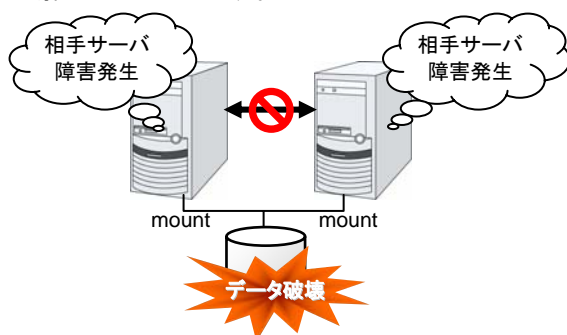


図 1-6 ネットワークパーティション症状

このような問題を「ネットワークパーティション症状」またはスプリットブレイン シンドローム (Split-brain-syndrome) と呼びます。フェイルオーバークラスタでは、すべてのインタコネクトが切断されたときに、確実に共有ディスク装置の排他制御を実現するためのさまざまな対応策が考えられています。

クラスタリソースの引き継ぎ

クラスタが管理するリソースにはディスク、IP アドレス、アプリケーションなどがあります。これらのクラスタリソースを引き継ぐための、フェイルオーバークラスタシステムの機能について説明します。

データの引き継ぎ

クラスタシステムでは、サーバ間で引き継ぐデータは共有ディスク装置上のパーティションに格納します。すなわち、データを引き継ぐとは、アプリケーションが使用するファイルが格納されているファイルシステムを健全なサーバ上でマウントしなおすことにほかなりません。共有ディスク装置は引き継ぐ先のサーバと物理的に接続されているので、クラスタソフトウェアが行うべきことはファイルシステムのマウントだけです。

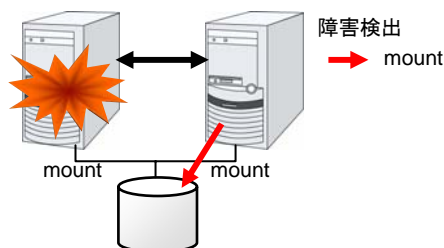


図 1-7 データの引き継ぎ

単純な話のようですが、クラスタシステムを設計・構築するうえで注意しなければならない点があります。

1 つは、ファイルシステムの復旧時間の問題です。引き継ごうとしているファイルシステムは、障害が発生する直前までほかのサーバで使用され、もしかしたらまさに更新中であつたかもしれません。このため、引き継ぐファイルシステムは通常ダーティであり、ファイルシステムの整合性チェックが必要な状態となっています。ファイルシステムのサイズが大きくなると、整合性チェックに必要な時間は莫大になり、場合によっては数時間もの時間がかかってしまいます。それがそのままフェイルオーバー時間(業務の引き継ぎ時間)に追加されてしまい、システムの可用性を低下させる要因になります。

もう 1 つは、書き込み保証の問題です。アプリケーションが大切なデータをファイルに書き込んだ場合、同期書き込みなどを利用してディスクへの書き込みを保証しようとします。ここでアプリケーションが書き込んだと思い込んだデータは、フェイルオーバー後にも引き継がれていることが期待されます。例えばメールサーバは、受信したメールをスプールに確実に書き込んだ時点で、クライアントまたはほかのメールサーバに受信完了を応答します。これによってサーバ障害発生後も、スプールされているメールをサーバの再起動後に再配信することができます。クラスタシステムでも同様に、一方のサーバがスプールへ書き込んだメールはフェイルオーバー後にもう一方のサーバが読み込めることを保証しなければなりません。

アプリケーションの引き継ぎ

クラスタソフトウェアが業務引き継ぎの最後に行う仕事は、アプリケーションの引き継ぎです。フォールトトレラントコンピュータ(FTC)とは異なり、一般的なフェイルオーバークラスタでは、アプリケーション実行中のメモリ内容を含むプロセス状態などを引き継ぎません。すなわち、障害が発生していたサーバで実行していたアプリケーションを健全なサーバで再実行することでアプリケーションの引き継ぎを行います。

例えば、データベース管理システム(DBMS)のインスタンスを引き継ぐ場合、インスタンスの起動時に自動的にデータベースの復旧(ロールフォワード/ロールバックなど)が行われます。このデータベース復旧に必要な時間は、DBMS のチェックポイントインターバルの設定などによってある程度の制御ができますが、一般的には数分程度必要となるようです。

多くのアプリケーションは再実行するだけで業務を再開できますが、障害発生後の業務復旧手順が必要なアプリケーションもあります。このようなアプリケーションのためにクラスタソフトウェアは業務復旧手順を記述できるよう、アプリケーションの起動の代わりにスクリプトを起動できるようにになっています。スクリプト内には、スクリプトの実行要因や実行サーバなどの情報をもとに、必要に応じて更新途中であつたファイルのクリーンアップなどの復旧手順を記述します。

フェイルオーバー総括

ここまでの内容から、次のようなクラスタソフトの動作が分かります。

- ◆ 障害検出(ハートビート/リソース監視)
- ◆ ネットワークパーティション症状解決(NP解決)¹
- ◆ クラスタ資源切り替え
 - データの引き継ぎ
 - IP アドレスの引き継ぎ
 - アプリケーションの引き継ぎ



図 1-8 フェイルオーバータイムチャート

クラスタソフトウェアは、フェイルオーバー実現のため、これらの様々な処置を 1 つ 1 つ確実に、短時間で実行することで、高可用性(High Availability)を実現しているのです。

注: Linux 版では図 1-8の「NP 解決」の時間はありません。

¹ Linux版では、この処理はありません。
セクション I CLUSTERPRO の概要

Single Point of Failure の排除

高可用性システムを構築するうえで、求められるもしくは目標とする可用性のレベルを把握することは重要です。これはすなわち、システムの稼働を阻害し得るさまざまな障害に対して、冗長構成をとることで稼働を継続したり、短い時間で稼働状態に復旧したりするなどの施策を費用対効果の面で検討し、システムを設計するということです。

Single Point of Failure(SPOF)とは、システム停止につながる部位を指す言葉であると前述しました。クラスタシステムではサーバの多重化を実現し、システムの SPOF を排除することができますが、共有ディスクなど、サーバ間で共有する部分については SPOF となり得ます。この共有部分を多重化もしくは排除するようシステム設計することが、高可用性システム構築の重要なポイントとなります。

クラスタシステムは可用性を向上させますが、フェイルオーバーには数分程度のシステム切り替え時間が必要となります。従って、フェイルオーバー時間は可用性の低下要因の 1 つともいえます。このため、高可用性システムでは、まず単体サーバの可用性を高める ECC メモリや冗長電源などの技術が本来重要なのですが、ここでは単体サーバの可用性向上技術には触れず、クラスタシステムにおいて SPOF となりがちな下記の 3 つについて掘り下げて、どのような対策があるか見ていきたいと思います。

- ◆ 共有ディスク
- ◆ 共有ディスクへのアクセスパス
- ◆ LAN

共有ディスク

通常、共有ディスクはディスクアレイにより RAID を組むので、ディスクのベアドライブは SPOF となりません。しかし、RAID コントローラを内蔵するため、コントローラが問題となります。多くのクラスタシステムで採用されている共有ディスクではコントローラの二重化が可能になっています。

二重化された RAID コントローラの利点を生かすためには、通常は共有ディスクへのアクセスパスの二重化を行う必要があります。ただし、二重化された複数のコントローラから同時に同一の論理ディスクユニット(LUN)へアクセスできるような共有ディスクの場合、それぞれのコントローラにサーバを 1 台ずつ接続すればコントローラ異常発生時にノード間フェイルオーバーを発生させることで高可用性を実現できます。

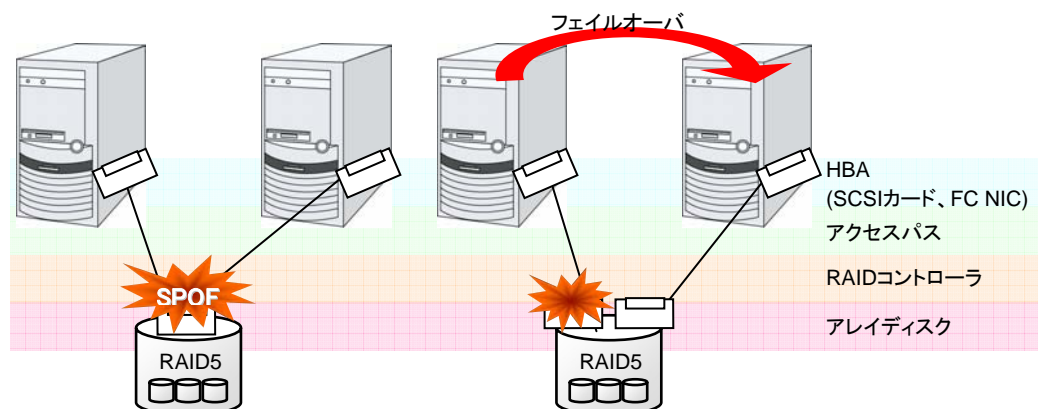


図 1-9 共有ディスクの RAID コントローラとアクセスパスが SPOF となっている例(左)と RAID コントローラとアクセスパスを分割した例

一方、共有ディスクを使用しないデータミラー型のフェイルオーバークラスタでは、すべてのデータをほかのサーバのディスクにミラーリングするため、SPOF が存在しない理想的なシステム構成を実現できます。ただし、欠点とはいえなくても、次のような点について考慮する必要があります。

- ◆ ネットワークを介してデータをミラーリングすることによるディスクI/O性能(特にwrite性能)
- ◆ サーバ障害後の復旧における、ミラー再同期中のシステム性能(ミラーコピーはバックグラウンドで実行される)
- ◆ ミラー再同期時間(ミラー再同期が完了するまでクラスタに組み込めない)

すなわち、データの参照が多く、データ容量が多くないシステムにおいては、データミラー型のフェイルオーバークラスタを採用するというのも可用性を向上させるポイントといえます。

共有ディスクへのアクセスパス

共有ディスク型クラスタの一般的な構成では、共有ディスクへのアクセスパスはクラスタを構成する各サーバで共有されます。SCSI を例に取れば、1 本の SCSI バス上に 2 台のサーバと共有ディスクを接続するということです。このため、共有ディスクへのアクセスパスの異常はシステム全体の停止要因となり得ます。

対策としては、共有ディスクへのアクセスパスを複数用意することで冗長構成とし、アプリケーションには共有ディスクへのアクセスパスが 1 本であるかのように見せることが考えられます。これを実現するデバイスドライバをパスフェイルオーバードライバなどと呼びます (パスフェイルオーバードライバは共有ディスクベンダーが開発してリリースするケースが多いのですが、Linux 版のパスフェイルオーバードライバは開発途上であったりしてリリースされていないようです。現時点では前述のとおり、共有ディスクのアレイコントローラごとにサーバを接続することで共有ディスクへのアクセスパスを分割する手法が Linux クラスタにおいては可用性確保のポイントとなります)。

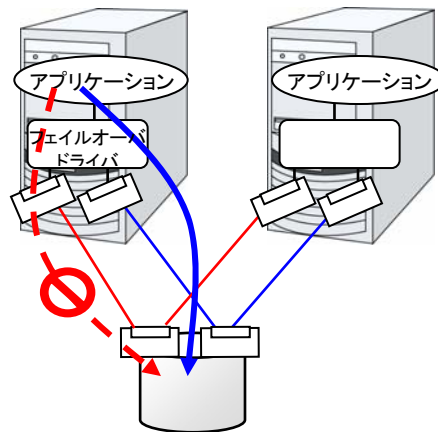


図 1-10 パスフェイルオーバードライバ

LAN

クラスタシステムに限らず、ネットワーク上で何らかのサービスを実行するシステムでは、LAN の障害はシステムの稼働を阻害する大きな要因です。クラスタシステムでは適切な設定を行えば NIC 障害時にノード間でフェイルオーバーを発生させて可用性を高めることは可能ですが、クラスタシステムの外側のネットワーク機器が故障した場合はやはりシステムの稼働を阻害します。

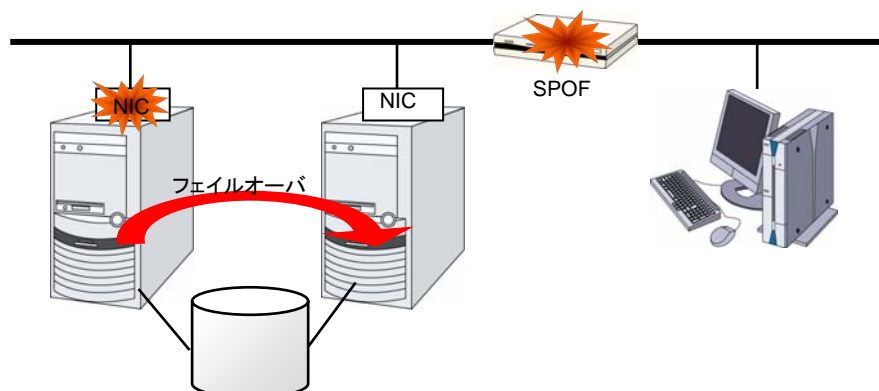


図 1-11 ルータが SPOF となる例

このようなケースでは、LAN を冗長化することでシステムの可用性を高めます。クラスタシステムにおいても、LAN の可用性向上には単体サーバでの技術がそのまま利用可能です。例えば、予備のネットワーク機器の電源を入れずに準備しておき、故障した場合に手動で入れ替えるといった原始的な手法や、高機能のネットワーク機器を冗長配置してネットワーク経路を多重化することで自動的に経路を切り替える方法が考えられます。また、インテル社の ANS ドライバのように NIC の冗長構成をサポートするドライバを利用することも考えられます。

ロードバランス装置 (Load Balance Appliance) やファイアウォールサーバ (Firewall Appliance) も SPOF となりやすいネットワーク機器です。これらもまた、標準もしくはオプションソフトウェアを利用することで、フェイルオーバー構成を組めるようになっているのが普通です。同時にこれらの機器は、システム全体の非常に重要な位置に存在するケースが多いため、冗長構成をとることはほぼ必須と考えるべきです。

可用性を支える運用

運用前評価

システムトラブルの発生要因の多くは、設定ミスや運用保守に起因するものであるともいわれています。このことから考えても、高可用性システムを実現するうえで運用前の評価と障害復旧マニュアルの整備はシステムの安定稼働にとって重要です。評価の観点としては、実運用に合わせて、次のようなことを実践することが可用性向上のポイントとなります。

- ◆ 障害発生箇所を洗い出し、対策を検討し、擬似障害評価を行い実証する
- ◆ クラスタのライフサイクルを想定した評価を行い、縮退運転時のパフォーマンスなどの検証を行う
- ◆ これらの評価をもとに、システム運用、障害復旧マニュアルを整備する

クラスタシステムの設計をシンプルにすることは、上記のような検証やマニュアルが単純化でき、システムの可用性向上のポイントとなることが分かります。

障害の監視

上記のような努力にもかかわらず障害は発生するものです。ハードウェアには経年劣化があり、ソフトウェアにはメモリリークなどの理由や設計当初のキャパシティプランニングを超えた運用をしてしまうことによる障害など、長期間運用を続ければ必ず障害が発生してしまいます。このため、ハードウェア、ソフトウェアの可用性向上と同時に、さらに重要となるのは障害を監視して障害発生時に適切に対処することです。万が一サーバに障害が発生した場合を例にとると、クラスタシステムを組むことで数分の切り替え時間でシステムの稼働を継続できますが、そのまま放置しておけばシステムは冗長性を失い次の障害発生時にはクラスタシステムは何の意味もなさなくなってしまう。

このため、障害が発生した場合、すぐさまシステム管理者は次の障害発生に備え、新たに発生した SPOF を取り除くなどの対処をしなければなりません。このようなシステム管理業務をサポートするうえで、リモートメンテナンスや障害の通報といった機能が重要になります。Linux では、リモートメンテナンスの面ではいうまでもなく非常に優れていますし、障害を通報する仕組みも整いつつあります。

以上、クラスタシステムを利用して高可用性を実現するうえで必要とされる周辺技術やそのほかのポイントについて説明しました。簡単にまとめると次のような点に注意しましょうということになるかと思います。

- ◆ Single Point of Failure を排除または把握する
- ◆ 障害に強いシンプルな設計を行い、運用前評価に基づき運用・障害復旧手順のマニュアルを整備する
- ◆ 発生した障害を早期に検出し適切に対処する

第 2 章 CLUSTERPRO の使用方法

本章では、CLUSTERPRO を構成するコンポーネントの説明と、クラスタシステムの設計から運用手順までの流れについて説明します。

本章で説明する項目は以下のとおりです。

• CLUSTERPRO とは?	18
• CLUSTERPRO の製品構成	18
• CLUSTERPRO のソフトウェア構成	18
• フェイルオーバーのしくみ	22
• リソースとは?	29
• CLUSTERPRO を始めよう!	31

CLUSTERPRO とは？

クラスタについて理解したところで、CLUSTERPRO の紹介を始めましょう。CLUSTERPRO とは、冗長化（クラスタ化）したシステム構成により、現用系のサーバでの障害が発生した場合に、自動的に待機系のサーバで業務を引き継がせることで、飛躍的にシステムの可用性と拡張性を高めることを可能にするソフトウェアです。

CLUSTERPRO の製品構成

CLUSTERPRO は大きく分けると 3 つのモジュールから構成されています。

- ◆ CLUSTERPRO Server

CLUSTERPRO の本体で、サーバの高可用性機能の全てが包含されています。また、WebManager のサーバ側機能も含まれます。

- ◆ CLUSTERPRO WebManager (WebManager)

CLUSTERPRO の運用管理を行うための管理ツールです。ユーザインターフェイスとして Web ブラウザを利用します。実体は CLUSTERPRO Server に組み込まれていますが、操作は管理端末上の Web ブラウザで行うため、CLUSTERPRO Server 本体とは区別されています。

- ◆ CLUSTERPRO Builder (Builder)

CLUSTERPRO の構成情報を作成するためのツールです。WebManager と同じく、ユーザインターフェイスとして Web ブラウザを利用します。Builder は Builder を利用する端末上で、CLUSTERPRO Server とは別にインストールする必要があります。

CLUSTERPRO のソフトウェア構成

CLUSTERPRO のソフトウェア構成は次の図のようになります。Linux サーバ上には「CLUSTERPRO Server (CLUSTERPRO 本体)」をインストールします。Builder は、管理 PC、あるいはサーバ上にインストールします。WebManager の本体 (CLUSTERPRO Server) は、CLUSTERPRO 本体と同時にインストールされているため、別途インストールする必要はありません。WebManager を利用する端末 (ブラウザ) は勿論管理 PC 上でも問題ありません。

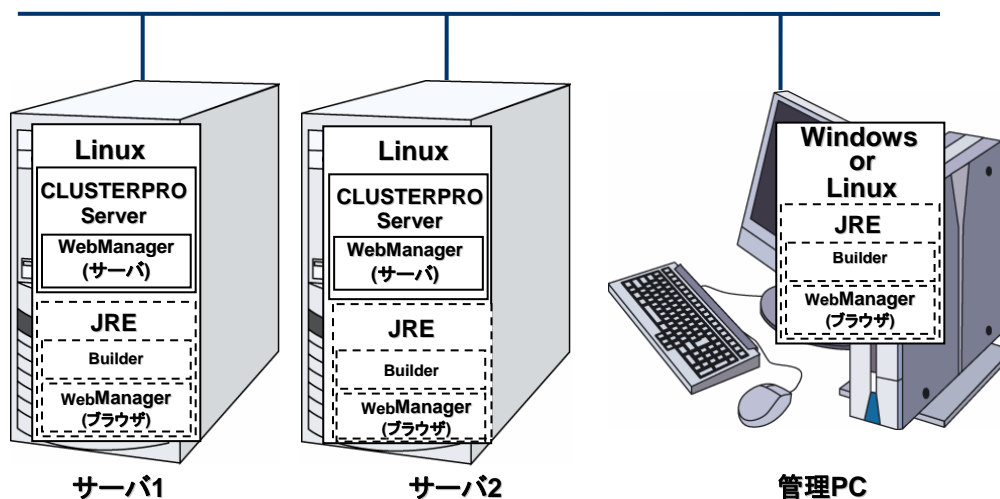


図 2-1 CLUSTERPRO のソフトウェア構成

CLUSTERPRO の障害監視のしくみ

CLUSTERPRO では、サーバ監視、業務監視、内部監視の 3 つの監視を行うことで、迅速かつ確実な障害検出を実現しています。以下にその監視の詳細を示します。

サーバ監視とは

サーバ監視とはフェイルオーバー型クラスタシステムの最も基本的な監視機能で、クラスタを構成するサーバが停止していないかを監視する機能です。

CLUSTERPRO はサーバ監視のために、定期的にサーバ同士で生存確認を行います。この生存確認をハートビートと呼びます。ハートビートは以下の通信パスを使用して行います。

◆ インタコネクト専用LAN

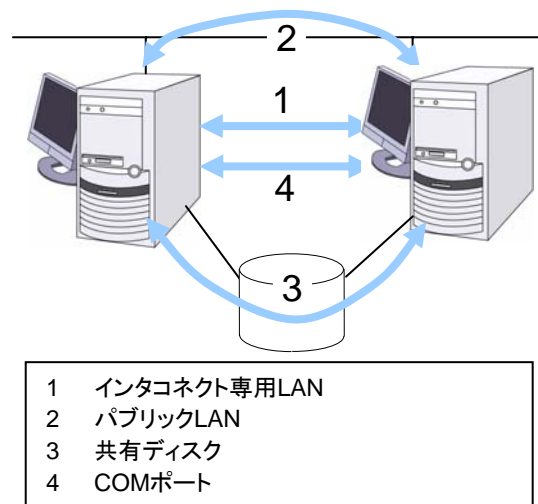
フェイルオーバー型クラスタ専用の通信パスで、一般の Ethernet NIC を使用します。ハートビートを行うと同時にサーバ間の情報交換に使用します。

◆ パブリックLAN

クライアントとの通信に使用している通信パスを予備のインタコネクトとして使用します。TCP/IP が使用できる NIC であればどのようなものでも構いません。インタコネクト専用 LAN インタコネクトサーバ間の情報交換にも使用します。

◆ 共有ディスク

フェイルオーバー型クラスタを構成する全てのサーバに接続されたディスク上に、CLUSTERPRO 専用のパーティション(CLUSTER パーティション)を作成し、CLUSTER パーティション上でハートビートを行います。



◆ COM ポート

フェイルオーバー型クラスタを構成するサーバ間を、COM ポートを介してハートビート通信を行い、他サーバの生存を確認します。

これらの通信経路を使用することでサーバ間の通信の信頼性は飛躍的に向上し、ネットワークパーティション症状の発生を防ぎます。

注：ネットワークパーティション症状(Split-brain-syndrome)について：クラスタサーバ間の全ての通信路に障害が発生しネットワーク的に分断されてしまう状態のことです。ネットワークパーティション症状に対応できていないクラスタシステムでは、通信路の障害とサーバの障害を区別できず、同一資源を複数のサーバからアクセスしデータ破壊を引き起こす場合があります。

業務監視とは

業務監視とは、業務アプリケーションそのものや業務が実行できない状態に陥る障害要因を監視する機能です。

◆ アプリケーションの死活監視

アプリケーションを起動用のリソース (EXEC リソースと呼びます) により起動を行い、監視用のリソース (PID モニタリソースと呼びます) により定期的にプロセスの生存を確認することで実現します。業務停止要因が業務アプリケーションの異常終了である場合に有効です。

注：

- CLUSTERPRO が直接起動したアプリケーションが監視対象の常駐プロセスを起動し終了してしまうようなアプリケーションでは、常駐プロセスの異常を検出することはできません。
 - アプリケーションの内部状態の異常 (アプリケーションのストールや結果異常) を検出することはできません。
-

◆ リソースの監視

CLUSTERPRO のモニタリソースによりクラスタリソース(ディスクパーティション、IP アドレスなど)やパブリック LAN の状態を監視することで実現します。業務停止要因が業務に必要なリソースの異常である場合に有効です。

内部監視とは

内部監視とは、CLUSTERPRO 内部のモジュール間相互監視です。CLUSTERPRO の各監視機能が正常に動作していることを監視します。

次のような監視を CLUSTERPRO 内部で行っています。

◆ CLUSTERPROプロセスの死活監視

監視できる障害と監視できない障害

CLUSTERPRO には、監視できる障害とできない障害があります。クラスタシステム構築時、運用時に、どのような監視が検出可能なのか、または検出できないのかを把握しておくことが重要です。

サーバ監視で検出できる障害とできない障害

監視条件: 障害サーバからのハートビートが途絶

- ◆ 監視できる障害の例
 - ハードウェア障害(OS が継続動作できないもの)
 - panic
- ◆ 監視できない障害の例
 - OS の部分的な機能障害(マウス/キーボードのみが動作しない等)

業務監視で検出できる障害とできない障害

監視条件: 障害アプリケーションの消滅、継続的なリソース異常、あるネットワーク装置への通信路切断

- ◆ 監視できる障害の例
 - アプリケーションの異常終了
 - 共有ディスクへのアクセス障害(HBA²の故障など)
 - パブリック LAN NIC の故障
- ◆ 監視できない障害の例
 - アプリケーションのストール/結果異常

アプリケーションのストール/結果異常を CLUSTERPRO で直接監視することはできませんが、アプリケーションを監視し異常検出時に自分自身を終了するプログラムを作成し、そのプログラムを EXEC リソースで起動、PID モニタリソースで監視することで、フェイルオーバを発生させることは可能です。

² Host Bus Adapterの略で、共有ディスク側ではなく、サーバ本体側のアダプタのことです。
セクション I CLUSTERPRO の概要

フェイルオーバーのしくみ

CLUSTERPRO は障害を検出すると、フェイルオーバー開始前に検出した障害がサーバの障害かネットワークパーティション症状かを判別します。この後、健全なサーバ上で各種リソースを活性化し業務アプリケーションを起動することでフェイルオーバーを実行します。

このとき、同時に移動するリソースの集まりをフェイルオーバーグループと呼びます。フェイルオーバーグループは利用者から見た場合、仮想的なコンピュータとみなすことができます。

注: クラスタシステムでは、アプリケーションを健全なノードで起動しなおすことでフェイルオーバーを実行します。このため、アプリケーションのメモリ上に格納されている実行状態をフェイルオーバーすることはできません。

障害発生からフェイルオーバー完了までの時間は数分間必要です。以下にタイムチャートを示します。

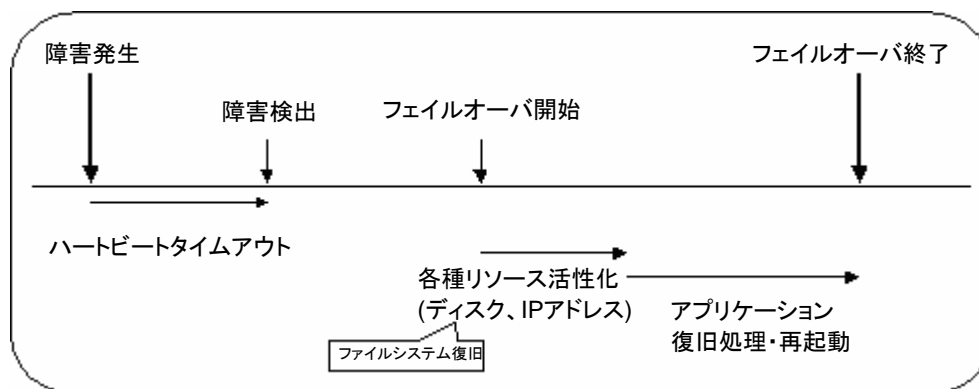


図 2-2 フェイルオーバーのタイムチャート

- ◆ **ハートビートタイムアウト**
 - ・ 業務を実行しているサーバの障害発生後、待機系がその障害を検出するまでの時間です。
 - ・ 業務の負荷に応じてクラスタプロパティの設定値を調整します。
(出荷時設定では 90 秒に設定されています。)
- ◆ **各種リソース活性化**
 - ・ 業務で必要なリソースを活性化するための時間です。
 - ・ 一般的な設定では数秒で活性化しますが、フェイルオーバーグループに登録されているリソースの種類や数によって必要時間は変化します。
(詳しくは、『CLUSTERPRO インストール & 設定ガイド』を参照してください。)
- ◆ **開始スクリプト実行時間**
 - ・ データベースのロールバック/ロールフォワードなどのデータ復旧時間と業務で使用するアプリケーションの起動時間です。
 - ・ ロールバック/ロールフォワード時間などはチェックポイントインターバルの調整である程度予測可能です。詳しくは、各ソフトウェア製品のドキュメントを参照してください。

フェイルオーバーリソース

CLUSTERPRO がフェイルオーバー対象とできる主なリソースは以下のとおりです。

- ◆ 切替パーティション (ディスクリソース、ミラーディスクリソースなど)
 - 業務アプリケーションが引き継ぐべきデータを格納するためのディスクパーティションです。
- ◆ フローティングIPアドレス (フローティングIPリソース)
 - フローティング IP アドレスを使用して業務へ接続することで、フェイルオーバーによる業務の実行位置(サーバ)の変化をクライアントは気にする必要がなくなります。
 - パブリック LAN アダプタへの IP アドレス動的割り当てと ARP パケットの送信により実現しています。ほとんどのネットワーク機器からフローティング IP アドレスによる接続が可能です
- ◆ スクリプト (EXEC リソース)
 - CLUSTERPRO では、業務アプリケーションをスクリプトから起動します。
 - 共有ディスクにて引き継がれたファイルはファイルシステムとして正常であっても、データとして不完全な状態にある場合があります。スクリプトにはアプリケーションの起動のほか、フェイルオーバー時の業務固有の復旧処理も記述します。

注：クラスタシステムでは、アプリケーションを健全なノードで起動しなおすことでフェイルオーバーを実行します。このため、アプリケーションのメモリ上に格納されている実行状態をフェイルオーバーすることはできません。

フェイルオーバー型クラスタのシステム構成

フェイルオーバー型クラスタは、ディスクアレイ装置をクラスタサーバ間で共有します。サーバ障害時には待機系サーバが共有ディスク上のデータを使用し業務を引き継ぎます。

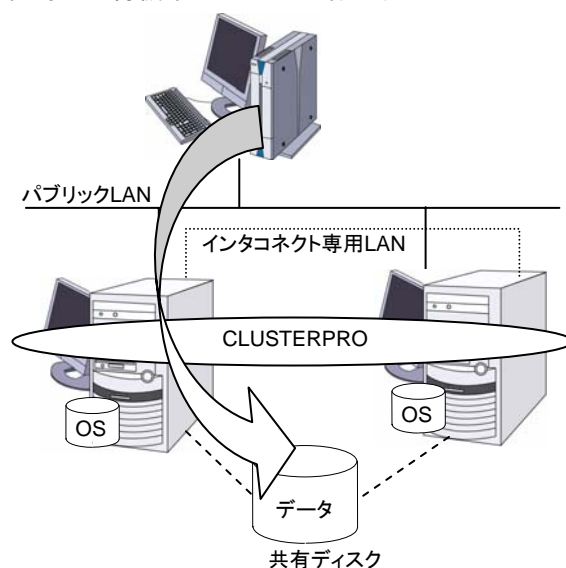


図 2-3 システム構成

セクション I CLUSTERPRO の概要

フェイルオーバー型クラスタでは、運用形態により、次のように分類できます。

片方向スタンバイクラスタ

一方のサーバを現用系として業務を稼動させ、他方のサーバを待機系として業務を稼動させない運用形態です。最もシンプルな運用形態でフェイルオーバー後の性能劣化のない可用性の高いシステムを構築できます。

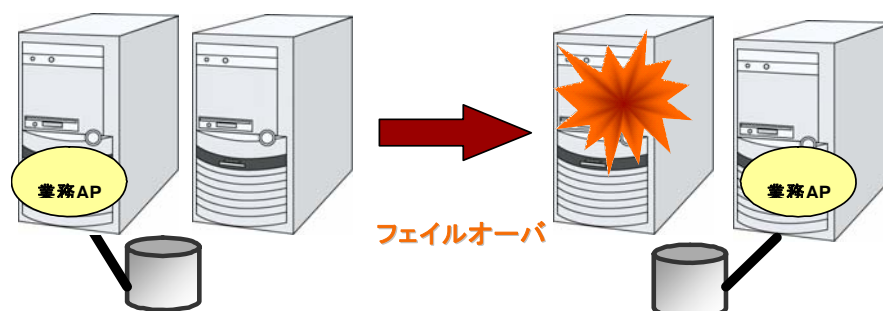
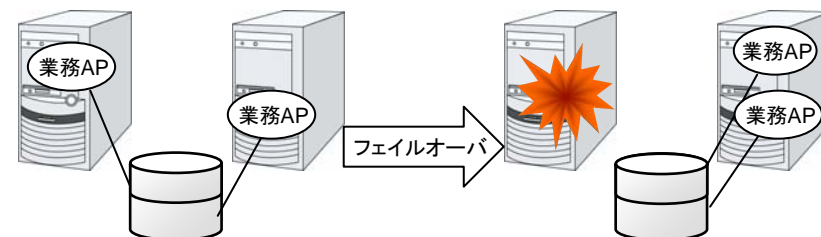


図 2-4 片方向スタンバイクラスタ

同一アプリケーション双方向スタンバイクラスタ

複数のサーバである業務アプリケーションを稼動させ相互に待機する運用形態です。アプリケーションは双方向スタンバイ運用をサポートしているものでなければなりません。ある業務データを複数に分割できる場合に、アクセスしようとしているデータによってクライアントからの接続先サーバを変更することで、データ分割単位での負荷分散システムを構築できます。

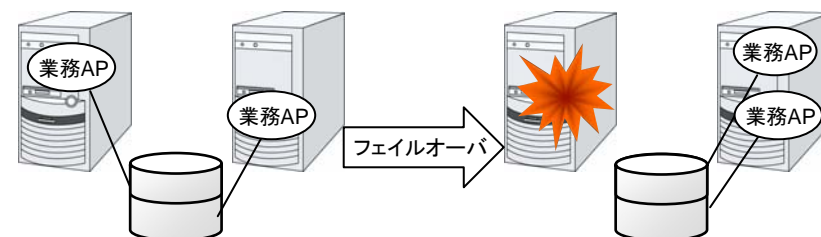


- ※ 図の業務APは同一アプリケーション
- ※ フェイルオーバー後にひとつのサーバ上で複数の業務APインスタンスが動く

図 2-5 同一アプリケーション双方向スタンバイクラスタ

異種アプリケーション双方向スタンバイクラスタ

複数の種類の業務アプリケーションをそれぞれ異なるサーバで稼動させ相互に待機する運用形態です。アプリケーションが双方向スタンバイ運用をサポートしている必要はありません。業務単位での負荷分散システムを構築できます。

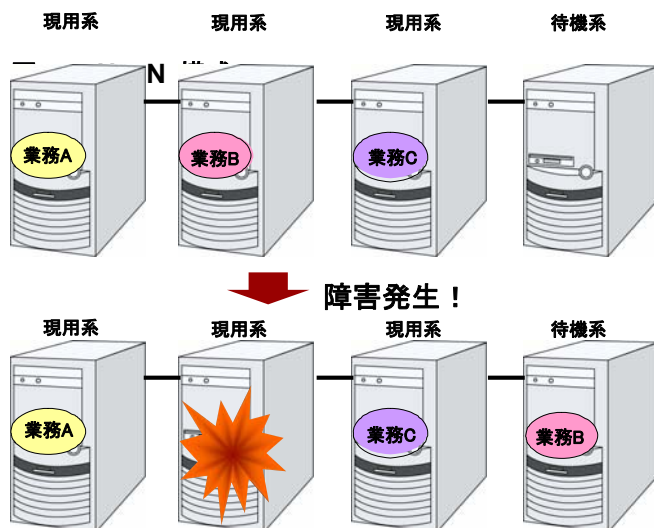


- ※ 業務1と業務2は異なるアプリケーションを使用

図 2-6 異種アプリケーション双方向スタンバイクラスタ

N + N 構成

ここまでの構成を応用し、より多くのノードを使用した構成に拡張することも可能です。下図は、3種の業務を3台のサーバで実行し、いざ問題が発生した時には1台の待機系にその業務を引き継ぐという構成です。片方向スタンバイでは、正常時のリソースの無駄は 1/2 でしたが、この構成なら正常時の無駄を 1/4 まで削減でき、かつ、1台までの異常発生であればパフォーマンスの低下もありません。



共有ディスク型のハードウェア構成

共有ディスク構成の CLUSTERPRO の HW 構成は下図のようになります。

サーバ間の通信用に

- ◆ NICを2枚 (1枚は外部との通信と流用、1枚はCLUSTERPRO専用)
- ◆ RS232Cクロスケーブルで接続されたCOMポート
- ◆ 共有ディスクの特定領域

を利用する構成が一般的です。

共有ディスクとの接続インターフェイスは SCSI か FibreChannel ですが、最近では FibreChannel による接続が一般的です。

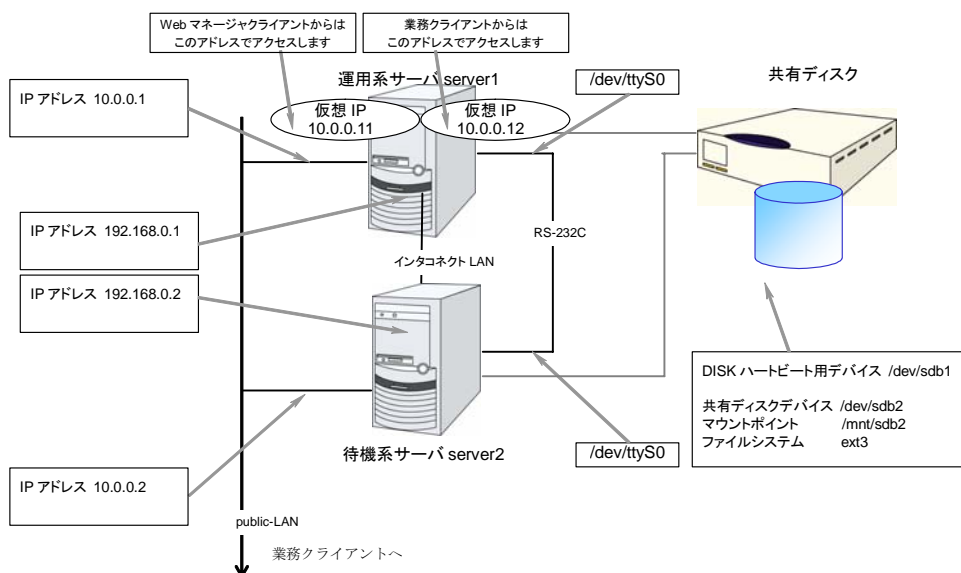


図 2-8 共有ディスク使用時のクラスタ環境のサンプル

ミラーディスク型のハードウェア構成

データミラー構成の CLUSTERPRO は、下図のような構成になります。

共有ディスク構成と比べ、ミラーディスクデータコピー用のネットワークが必要となりますが、通常、CLUSTERPRO の内部通信用 NIC と兼用します。

また、ミラーディスクは接続インターフェイス(IDE or SCSI)には依存しません。

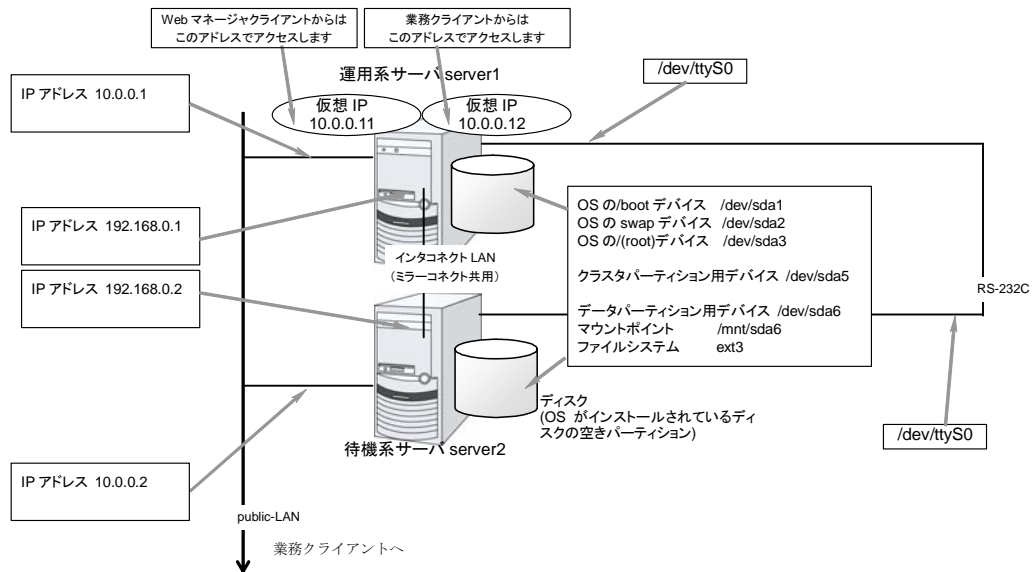


図 2-9 ミラーディスク使用時のクラスタ環境のサンプル(OS がインストールされているディスクにクラスタパーティション、データパーティションを確保する場合)

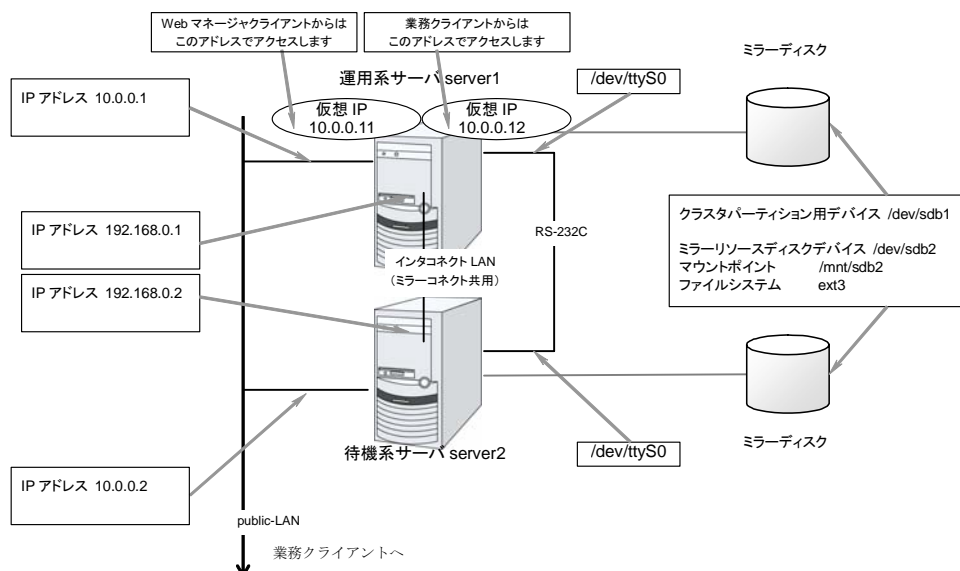


図 2-10 ミラーディスク使用時のクラスタ環境のサンプル(クラスタパーティション、データパーティション用のディスクを用意する場合)

クラスタオブジェクトとは？

CLUSTERPRO では各種リソースを下のような構成で管理しています。

- ◆ クラスタオブジェクト
クラスタの構成単位となります。
- ◆ サーバオブジェクト
実体サーバを示すオブジェクトで、クラスタオブジェクトに属します。
- ◆ ハートビートリソースオブジェクト
実体サーバのNW部分を示すオブジェクトで、サーバオブジェクトに属します。
- ◆ グループオブジェクト
仮想サーバを示すオブジェクトで、クラスタオブジェクトに属します。
- ◆ グループリソースオブジェクト
仮想サーバの持つリソース(NW、ディスク)を示すオブジェクトでグループオブジェクトに属します。
- ◆ モニタリソースオブジェクト
監視機構を示すオブジェクトで、クラスタオブジェクトに属します。

リソースとは?

CLUSTERPRO では、監視する側とされる側の対象をすべてリソースと呼び、分類して管理します。このことにより、より明確に監視/被監視の対象を区別できるほか、クラスタ構築や障害検出時の対応が容易になります。リソースはハートビートリソース、グループリソース、モニタリソースの 3 つに分類されます。以下にその概略を示します。

ハートビートリソース

サーバ間で、お互いの生存を確認するためのリソースです。

以下に現在サポートされているハートビートリソースを示します。

- ◆ LANハートビートリソース
Ethernetを利用した通信を示します。
- ◆ カーネルモードLANハートビートリソース
Ethernetを利用した通信を示します。
- ◆ COMハートビートリソース
RS232C(COM)を利用した通信を示します。
- ◆ ディスクハートビートリソース
共有ディスク上の特定パーティション(ディスクハートビート用パーティション)を利用した通信を示します。共有ディスク構成の場合のみ利用可能です。

グループリソース

フェイルオーバーを行う際の単位となる、フェイルオーバーグループを構成するリソースです。

以下に現在サポートされているグループリソースを示します。

- ◆ フローティングIPリソース (fip)
仮想的なIPアドレスを提供します。クライアントからは一般のIPアドレスと同様にアクセス可能です。
- ◆ EXECリソース (exec)
業務(DB、httpd、etc..)を起動／停止するための仕組みを提供します。
- ◆ ディスクリソース (disk)
共有ディスク上の指定パーティションを提供します。(共有ディスク)構成の場合のみ利用可能です。
- ◆ ミラーディスクリソース (md)
ミラーディスク上の指定パーティションを提供します。(ミラーディスク)構成の場合のみ利用可能です。
- ◆ RAWリソース (raw)
共有ディスク上のRAWデバイスを提供します。共有ディスク構成の場合のみ利用可能です。
- ◆ VxVMディスクグループリソース (vxdg)
共有ディスク上のVxVMディスクグループを提供します。VxVMボリュームリソースと共に使用します。(共有ディスク)構成の場合のみ利用可能です。

- ◆ VxVMボリュームリソース (vxvol)
共有ディスク上のVxVMボリュームを提供します。VxVMディスクグループリソースと共に使用します。(共有ディスク)構成の場合のみ利用可能です。
- ◆ NASリソース (nas)
NASサーバ上の共有リソースへ接続します。(クラスタサーバがNASのサーバ側として振る舞うリソースではありません。)

モニタリソース

クラスタシステム内で、監視を行う主体であるリソースです。

以下に現在サポートされているモニタリソースを示します。

- ◆ IPモニタリソース (ipw)
外部のIPアドレスの監視機構を提供します。
- ◆ ディスクモニタリソース (diskw)
ディスクの監視機構を提供します。共有ディスクの監視にも利用されます。
- ◆ ミラーディスクモニタリソース (mdw)
ミラーディスクの監視機構を提供します。
- ◆ ミラーディスクコネクタモニタリソース (mdnw)
ミラーディスクコネクタの監視機構を提供します。
- ◆ PIDモニタリソース (pidw)
EXECリソースで起動したプロセスの死活監視機能を提供します。
- ◆ ユーザ空間モニタリソース (userw)
ユーザ空間のストール監視機構を提供します。
- ◆ RAWモニタリソース (raww)
ディスクの監視機構を提供します。RAWデバイスを使用するためreadサイズが小さいのでシステムへの負荷が軽減できます。共有ディスクの監視にも利用されます。
- ◆ NIC Link Up/Downモニタリソース (miiw)
LANケーブルのリンクステータスの監視機構を提供します。
- ◆ VxVMデーモンモニタリソース (vxdw)
VxVMのデーモンの監視機構を提供します。(共有ディスク)構成の場合のみ利用可能です。
- ◆ VxVMボリュームモニタリソース (vxvolw)
VxVMのボリュームの監視機構を提供します。(共有ディスク)構成の場合のみ利用可能です。
- ◆ マルチターゲットモニタリソース (mtw)
複数のモニタリソースを束ねたステータスを提供します。

CLUSTERPRO を始めよう!

以上で CLUSTERPRO の簡単な説明が終了しました。

以降は、以下の流れに従い、対応するガイドを読み進めながら CLUSTERPRO を使用したクラスタシステムの構築を行ってください。

最新情報の確認

本ガイドのセクション II 『リリースノート (CLUSTERPRO 最新情報)』を参照してください。

クラスタシステムの設計

『インストール&設定ガイド』の「セクション I クラスタシステムの設計」および『リファレンスガイド』の「セクション II リソース詳細」を参照してください。

クラスタシステムの構築

『インストール&設定ガイド』の全編を参照してください。

オプションの監視コマンドを使用する場合は、監視対象アプリケーション別の『管理者ガイド』を参照してください。

クラスタシステムの運用開始後の障害対応

『リファレンスガイド』の「セクション III メンテナンス情報」を参照してください。

セクション II リリースノート (CLUSTERPRO 最新情報)

このセクションでは、CLUSTERPRO の最新情報を記載します。サポートするハードウェアやソフトウェアについての最新の詳細情報を記載します。また、制限事項や、既知の問題とその回避策についても説明します。

- 第 3 章 CLUSTERPRO の動作環境
- 第 4 章 最新バージョン情報
- 第 5 章 注意制限事項
- 第 6 章 アップデート手順

第 3 章 CLUSTERPRO の動作環境

本章では、CLUSTERPRO の動作環境について説明します。

本章で説明する項目は以下の通りです。

• ハードウェア	36
• CLUSTERPRO Server の動作環境	38
• Builderの動作環境	42
• WebManagerの動作環境.....	43

ハードウェア

CLUSTERPRO は以下のアーキテクチャのサーバで動作します。

- ◆ IA32
- ◆ x86_64
- ◆ IA64 (Replicator, Agent, Alert Serviceは未サポート)
- ◆ ppc64 (Replicator, Agent, Alert Serviceは未サポート)

スペック

CLUSTERPRO Server で必要なスペックは下記の通りです。

- ◆ RS-232Cポート 1つ (3ノード以上のクラスタを構築する場合は不要)
- ◆ Ethernetポート 2つ以上
- ◆ 共有ディスク (Replicatorを使用する場合は不要)
- ◆ ミラー用ディスク または ミラー用空きパーティション (Replicatorを使用する場合は必要)
- ◆ CD-ROMドライブ

構築、構成変更時には Builder との情報のやりとりのため以下が必要です。

- ◆ FDドライブ、USBメモリなどのリムーバブルメディア
- ◆ Builderを動作させるマシンとファイルを共有する手段

動作確認済ディスクインターフェイス

Replicator のミラーディスクとして確認済みのディスクタイプは下記の通りです。

ディスクのタイプ	ホスト側ドライバ呼称	備考
IDE	ide	～120GBまで確認済
SCSI	aic7xxx	
SCSI	aic79xx	
SCSI	sym53c8xx	
SCSI	mptbase, mptscsih	
RAID	megaraid(SCSIタイプ)	
RAID	megaraid (IDEタイプ)	～275GBまで確認済
S-ATA	sata-nv	～80GBまで確認済
S-ATA	ata-piix	～120GBまで確認済

動作確認済ネットワークインターフェイス

Replicator のミラーディスクのミラーディスクコネクタ(ミラー通信で使用する系)として確認済みのネットワークボードは下記の通りです。

チップ呼称	ドライバ呼称
Intel 82557/8/9	e100
Intel 82540EM Intel 82544EI Intel 82546EB Intel 82546GB	e1000
Broadcom BCM5701 Broadcom BCM5703 Broadcom BCM5721	bcm5700

ここに掲載しているものは代表的な一例であり、これ以外の製品も利用可能です。

ソフトウェア

CLUSTERPRO Serverの動作環境

動作可能なディストリビューションとkernel

CLUSTERPRO 独自の kernel モジュールがあるため、CLUSTERPRO Serverの動作環境は kernel モジュールのバージョンに依存します。適合する kernel モジュール(ドライバ)を提供している kernel バージョンの情報を提示します。

下記以外のバージョンでは正常に動作しません。

IA32

ディストリビューション	kernel バージョン	Replicator サポート	clpka,clpkhb サポート	CLUSTERPRO Version	備考
Turbolinux 10 Server サービスパック	2.6.8-6 2.6.8-6smp 2.6.8-6smp64G	○	○	1.0.0-1~	
	2.6.8-12 2.6.8-12smp 2.6.8-12smp64G	○	○	1.0.2-1~	
Turbo ApplianceServer2.0	2.6.8-6 2.6.8-6smp 2.6.8-6smp64G	○	○	1.0.0-1~	
Red Hat Enterprise Linux AS/ES 4 (update3)	2.6.9-34.EL 2.6.9-34.ELsmp 2.6.9-34.ELhugemem	○	○	1.0.0-1~	1
Red Hat Enterprise Linux AS/ES 4 (update4)	2.6.9-42.EL 2.6.9-42.ELsmp 2.6.9-42.ELhugemem	○	○	1.0.1-1~	1
	2.6.9-42.0.3.EL 2.6.9-42.0.3.ELsmp 2.6.9-42.0.3.ELhugemem	○	○	1.0.2-1~	1
	2.6.9-42.0.8.EL 2.6.9-42.0.8.ELsmp 2.6.9-42.0.8.ELhugemem	○	○	1.0.2-1~	1
	2.6.9-42.0.10.EL 2.6.9-42.0.10.ELsmp 2.6.9-42.0.10.ELhugemem	○	○	1.0.2-1~	1
MIRACLE LINUX V4.0	2.6.9-11.25AX 2.6.9-11.25AXsmp 2.6.9-11.25AXhugemem	○	○	1.0.0-1~	
MIRACLE LINUX V4.0 SP1	2.6.9-34.21AX 2.6.9-34.21AXsmp 2.6.9-34.21AXhugemem	○	○	1.0.0-1~	
	2.6.9-34.28AX 2.6.9-34.28AXsmp 2.6.9-34.28AXhugemem	○	○	1.0.1-1~	
MIRACLE LINUX V4.0 SP2	2.6.9-42.7AX 2.6.9-42.7AXsmp 2.6.9-42.7AXhugemem	○	○	1.0.2-1~	

ディストリビューション	kernel バージョン	Replicator サポート	clpka,clpkhb サポート	CLUSTERPRO Version	備考
Novell SUSE LINUX Enterprise Server 9 (SP3)	2.6.5-7.244-default 2.6.5-7.244-smp 2.6.5-7.244-bigsmpt	○	○	1.0.0-1~	

x86_64

ディストリビューション	kernel バージョン	Replicator サポート	clpka,clpkhb サポート	CLUSTERPRO Version	備考
Turbolinux 10 Server サービスパック	2.6.13-8 2.6.13-8smp	○	○	1.0.0-1~	
	2.6.13-17 2.6.13-17smp	○	○	1.0.2-1~	
Red Hat Enterprise Linux AS/ES 4 (update3)	2.6.9-34.EL 2.6.9-34.ELsmp	○	○	1.0.0-1~	1
Red Hat Enterprise Linux AS/ES 4 (update4)	2.6.9-42.EL 2.6.9-42.ELsmp	○	○	1.0.1-1~	1
	2.6.9-42.0.3.EL 2.6.9-42.0.3.ELsmp	○	○	1.0.2-1~	1
	2.6.9-42.0.8.EL 2.6.9-42.0.8.ELsmp	○	○	1.0.2-1~	1
	2.6.9-42.0.10.EL 2.6.9-42.0.10.ELsmp	○	○	1.0.2-1~	1
MIRACLE LINUX V4.0	2.6.9-11.25AX 2.6.9-11.25AXsmp	○	○	1.0.0-1~	
MIRACLE LINUX V4.0 SP1	2.6.9-34.21AX 2.6.9-34.21AXsmp	○	○	1.0.0-1~	
	2.6.9-34.28AX 2.6.9-34.28AXsmp 2.6.9-34.28AXlargesmp	○	○	1.0.1-1~	
MIRACLE LINUX V4.0 SP2	2.6.9-42.7AX 2.6.9-42.7AXsmp	○	○	1.0.2-1~	
Novell SUSE LINUX Enterprise Server 9 (SP3)	2.6.5-7.244-default 2.6.5-7.244-smp	○	○	1.0.0-1~	

IA64

ディストリビューション	kernel バージョン	Replicator サポート	clpka,clpkhb サポート	CLUSTERPRO Version	備考
Red Hat Enterprise Linux AS/ES 4 (update3)	2.6.9-34.EL	×	○	1.0.0-1~	
Red Hat Enterprise Linux AS/ES 4 (update4)	2.6.9-42.EL	×	○	1.0.1-1~	
	2.6.9-42.0.3.EL	×	○	1.0.2-1~	
	2.6.9-42.0.8.EL	×	○	1.0.2-1~	
	2.6.9-42.0.10.EL	×	○	1.0.2-1~	
Asianux 2.0 SP1準拠 ディストリビューション	2.6.9-34.21AX	×	○	1.0.0-1~	
Novell SUSE LINUX Enterprise Server 9 (SP3)	2.6.5-7.244-default	×	○	1.0.0-1~	

ppc64

ディストリビューション	kernel バージョン	Replicator サポート	clpka,clpkhb 動作可否	CLUSTERPRO Version	備考
Red Hat Enterprise Linux AS/ES 4 (update3)	2.6.9-34.EL	×	○	1.0.0-1~	
Red Hat Enterprise Linux AS/ES 4 (update4)	2.6.9-42.EL	×	○	1.0.1-1~	
	2.6.9-42.0.3.EL	×	○	1.0.2-1~	
	2.6.9-42.0.8.EL	×	○	1.0.2-1~	
	2.6.9-42.0.10.EL	×	○	1.0.2-1~	
Asianux 2.0 SP1準拠 ディストリビューション	2.6.9-34.21AX	×	○	1.0.0-1~	
Novell SUSE LINUX Enterprise Server 9 (SP3)	2.6.5-7.244-default	×	○	1.0.0-1~	

備考欄凡例

- 1 CLUSTERPROのミラードライバがファイルシステム vxfs に対応しているkernelバージョンです。対応しているCLUSTERPROのバージョンは1.0.1-1以降です。CLUSTERPROのバージョン1.0.0-1ではファイルシステム vxfs に対応していません。

必要メモリ容量とディスクサイズ

	必要メモリサイズ		必要ディスクサイズ		備考
	ユーザモード	kernel モード	インストール直後	運用時最大	
IA32	35MB	32MB + 2MB × ミラーリソース数	32MB	600MB	
x86_64	40MB	32MB + 2MB × ミラーリソース数	32MB	600MB	
IA64	80MB	-	16MB	400MB	
ppc64	64MB	-	12MB	400MB	

Builder の動作環境

動作確認済OS、ブラウザ

最新情報は CLUSTERPRO のホームページで公開されている最新ドキュメントを参照してください。現在の対応状況は下記の通りです。

OS	ブラウザ	言語
Microsoft Windows® XP SP2 (IA32)	IE6 SP2	日本語/英語
Microsoft Windows Server™ 2003 SP1 以降(IA32)	IE6 SP1	日本語/英語
Novell SUSE LINUX Enterprise Server 9 SP2 (IA32)	FireFox 1.0.6	日本語/英語
Red Hat Enterprise Linux AS/ES 4 update3 (IA32)	FireFox 1.0.7	日本語/英語

注: 64bit、x86_64、ppc64 上では「Builder」は動作しません。構築時、構成変更時には 32bit マシンを用意してください。

Java実行環境

Builder を使用する場合には、Java 実行環境が必要です。

Sun Microsystems

Java(TM) Runtime Environment

Version 5.0 Update6 (1.5.0_06)以降

必要メモリ容量/ディスク容量

必要メモリ容量 32MB 以上

必要ディスク容量 2MB(Java 実行環境に必要な容量を除く)

対応するCLUSTERPROのバージョン

Builderバージョン	CLUSTERPRO X rpmバージョン
1.0.0-1	1.0.0-1
1.0.1-1	1.0.1-1 1.0.2-1

注: Builder のバージョンと CLUSTERPRO rpm バージョンは上記の対応表の組み合わせで使用してください。それ以外の組み合わせで使用すると正常に動作しない可能性があります。

WebManager の動作環境

動作確認済OS、ブラウザ

現在の対応状況は下記の通りです。

OS	ブラウザ	言語
Microsoft Windows® XP(IA32)	IE6 SP2	日本語/英語
Microsoft Windows Server™ 2003 SP1 以降(IA32, x86_64)	IE6 SP1	日本語/英語
Novell SUSE LINUX Enterprise Server 9 SP2 (IA32)	Firefox 1.0.6	日本語/英語
Red Hat Enterprise Linux AS/ES 4 update3 (IA32)	Firefox 1.0.7	日本語/英語

注: WebManager は 64bit、x86_64、ppc64 の Linux 上では 動作しません。Linux マシンでクラスタの管理をするには 32bit OS を用意してください。

Java実行環境

WebManager を使用する場合には、Java 実行環境が必要です。

Sun Microsystems

Java(TM) Runtime Environment

Version 5.0 Update6 (1.5.0_06)以降

必要メモリ容量/ディスク容量

必要メモリ容量 40MB 以上

必要ディスク容量 300KB(Java 実行環境に必要な容量を除く)

第 4 章 最新バージョン情報

本章では、CLUSTERPRO の最新情報について説明します。新しいリリースで強化された点、改善された点などをご紹介します。

現在は初期バージョンのため、ありません。

第 5 章 注意制限事項

本章では、注意事項や既知の問題とその回避策について説明します。

本章で説明する項目は以下の通りです。

• システム構成検討時	48
• OS インストール前、OS インストール時	51
• OS インストール後、CLUSTERPRO インストール前	54
• CLUSTERPRO の情報作成時	61
• CLUSTERPRO 運用後	64

システム構成検討時の注意事項

HW の手配、システム構成、共有ディスクの構成時に留意すべき事項について説明します。

Builder、WebManagerの動作OSについて

- ◆ 64bit、x86_64、ppc64上では「Builder」は動作しません。構築時、構成変更時には32bitマシンを用意してください。
- ◆ WebManager は64bit、x86_64、ppc64 のLinux上では 動作しません。Linuxマシンでクラスタの管理をするには32bit OSを用意してください。

ミラーディスクの要件について

- ◆ ミラーリソースとして使用するディスクはLinuxのmdやLVMによるストライプセット、ボリュームセット、ミラーリング、パリティ付ストライプセットの機能はサポートしていません。
- ◆ ミラーリソースを使用するにはミラー用のパーティション(データパーティションとクラスタパーティション)が必要です。
- ◆ ミラー用のパーティションの確保の方法は以下の2つがあります。
 - OS(rootパーティションやswapパーティション)と同じディスク上にミラー用のパーティション(クラスタパーティションとデータパーティション)を確保する
 - OS とは別のディスク(または LUN)を用意(追加)してミラー用のパーティションを確保する
- ◆ 以下を参考に上記を選定してください。
 - 障害時の保守性、性能を重視する場合
- OS とは別にミラー用のディスクを用意することを推奨します。
 - H/W Raid の仕様の制限で LUN の追加ができない場合
H/W Raid のブラインストールモデルで LUN 構成変更が困難な場合
- OS と同じディスクにミラー用のパーティションを確保します。
- ◆ ミラーリソースを複数使用する場合には、さらにミラーリソース毎に個別のディスクを用意(追加)することを推奨します。
同一のディスク上に複数のミラーリソースを確保すると性能の低下やミラー復旧に時間がかかることがあります。これらの現象はLinux OSのディスクアクセスの性能に起因するものです。
- ◆ ミラー用のディスクとして使用するにはディスクをサーバ間で同じにする必要があります。
 - ディスクのタイプ
両サーバのミラーディスクまたは、ミラー用のパーティションを確保するディスクは、ディスクのタイプを同じにしてください。
動作確認済みのディスクのタイプについては 36 ページの「動作確認済ディスクインターフェイス」を参照してください。

例)

組み合わせ	サーバ1	サーバ2
OK	SCSI	SCSI
OK	IDE	IDE
NG	IDE	SCSI

◆ ミラー用のディスクとして使用するディスクのジオメトリがサーバ間で異なる場合の注意

- ディスクのジオメトリ

両サーバのミラーディスクまたは、ミラー用のパーティションを確保するディスクは、ディスクのジオメトリが等しいものを推奨します。

- ジオメトリが異なる場合の注意事項

fdisk コマンドなどで確保したパーティションサイズはシリンダあたりのブロック(ユニット)数でアラインされます。

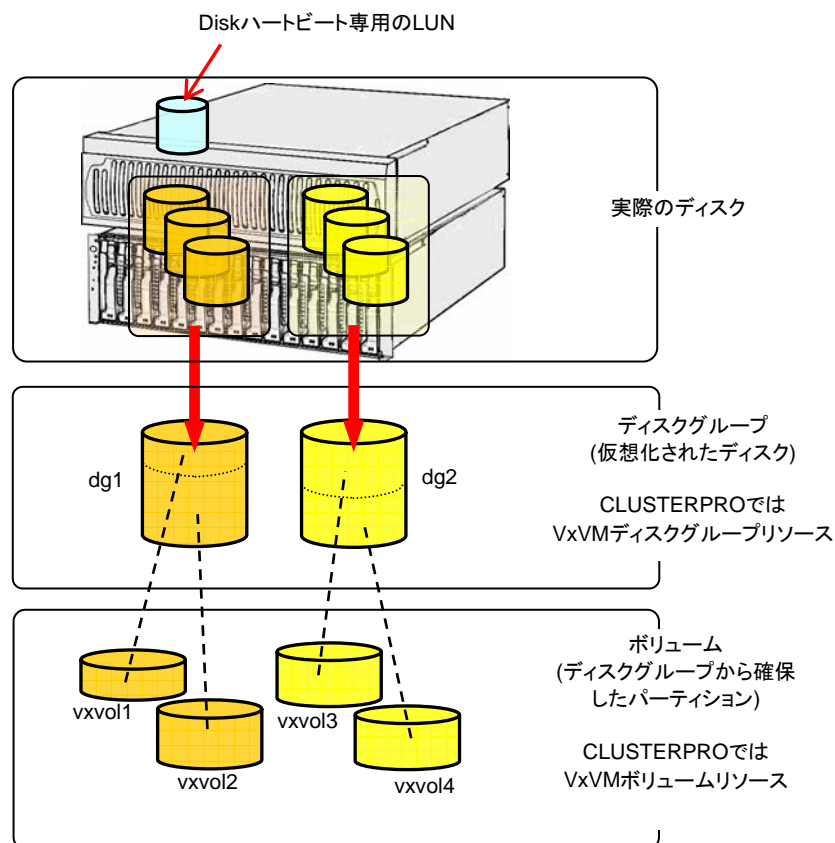
データパーティションのサイズと初期ミラー構築の方向の関係が以下になるようにデータパーティションを確保してください。

コピー元のサーバ ≤ コピー先のサーバ

コピー元のサーバとは、ミラーリソースが所属するフェイルオーバーグループのフェイルオーバーポリシーが高いサーバを指します。コピー先のサーバとは、ミラーリソースが所属するフェイルオーバーグループのフェイルオーバーポリシーが低いサーバを指します。

共有ディスクの要件について

- ◆ 共有ディスクはLinuxのmdlによるストライプセット、ボリュームセット、ミラーリング、パリティ付ストライプセットの機能はサポートしていません。
- ◆ 共有ディスクでLinuxのLVMによるストライプセット、ボリュームセット、ミラーリング、パリティ付ストライプセットの機能を使用する場合、以下の制限があります。
 - ディスクリソースに設定されたパーティションの ReadOnly,ReadWrite の制御を CLUSTERPRO が行うことができません。
ディスクリソースが活性化されていないサーバで誤って同じファイルシステムをマウントしないように運用回避をしてください。
 - LVM の論理ボリュームをディスクモニタリソースで監視できません。
ディスクモニタリソースの監視先として LVM のボリュームを構成する実デバイスを設定し、マルチターゲットモニタを組み合わせで監視をしてください。
- ◆ VxVMを使用する場合、CLUSTERPROのディスクハートビート用に共有ディスク上に、VxVMで制御対象としないLUNが必要です。共有ディスクのLUNの設計時に留意してください。



NIC Link Up/Downモニタリソース

NIC のボード、ドライバによっては、必要な `ioctl()` がサポートされていない場合があります。
 その場合には このモニタリソースは使用できません。

ミラーリソースのwrite性能について

- ◆ ミラーディスクのwrite処理はネットワークを経由して相手サーバのディスクへwrite、自サーバのディスクへwriteを行います。readは自サーバ側のディスクからのみreadします。
- ◆ 上記の理由により、クラスタ化していない単体サーバと比べてwrite性能が劣化します。
 writeに単体サーバ並みに高スループットが要求されるシステム(更新系が多いデータベースシステムなど)には共有ディスク使用をご提案ください。

OS インストール前、OS インストール時

OS をインストールするときに決定するパラメータ、リソースの確保、ネーミングルールなどで留意して頂きたいことです。

/opt/nec/clusterproのファイルシステムについて

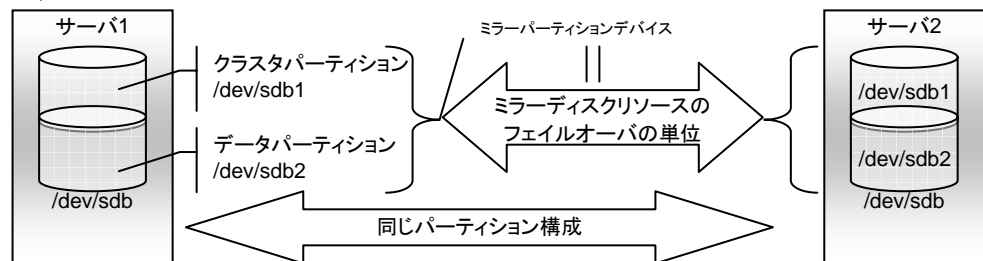
システムの対障害性の向上のために、ジャーナル機能を持つファイルシステムを使用することを推奨します。

ミラー用のディスクについて

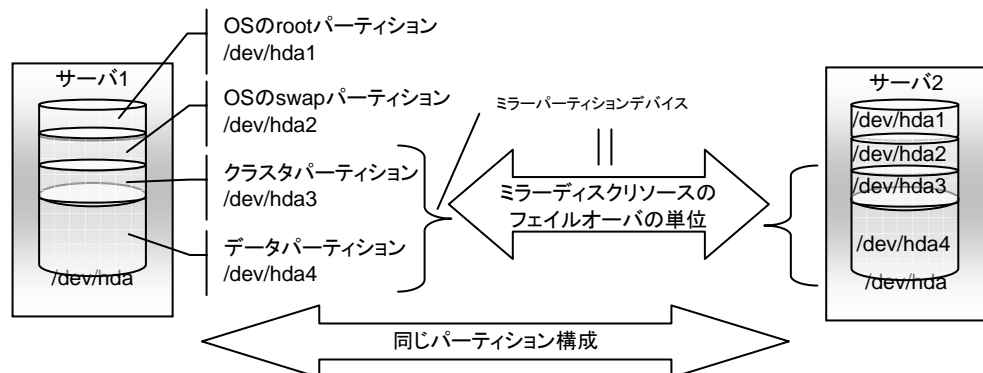
◆ ディスクのパーティション

両サーバで同一パーティションに対して、同一デバイス名でアクセスできるように設定してください。

例)両サーバに1つの SCSI ディスクを増設して1つのミラーディスクのペアにする場合



例)両サーバの OS が格納されている IDE ディスクの空き領域を使用してミラーディスクのペアにする場合



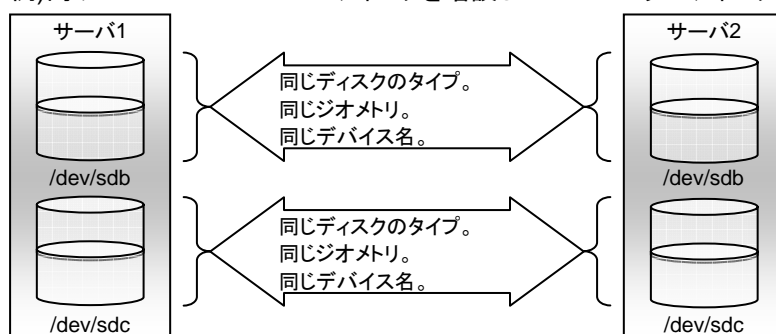
- ミラーパーティションデバイスは CLUSTERPRO のミラーリングドライバが上位に提供するデバイスです。
- クラスタパーティションとデータパーティションの 2 つのパーティションをペアで確保してください。
- OS(root パーティションや swap パーティション)と同じディスク上にミラーパーティション(クラスタパーティション、データパーティション)を確保することも可能です。
 - 障害時の保守性、性能を重視する場合
OS(root パーティションや swap パーティション)と別にミラー用のディスクを用意することを推奨します。
 - H/W Raid の仕様の制限で LUN の追加ができない場合
H/W Raid のブリーインストールモデルで LUN 構成変更が困難な場合
OS(root パーティションや swap パーティション)と同じディスクにミラーパーティション(クラスタパーティション、データパーティション)を確保することも可能です。

◆ ディスクの配置

ミラーディスクとして複数のディスクを使用することができます。

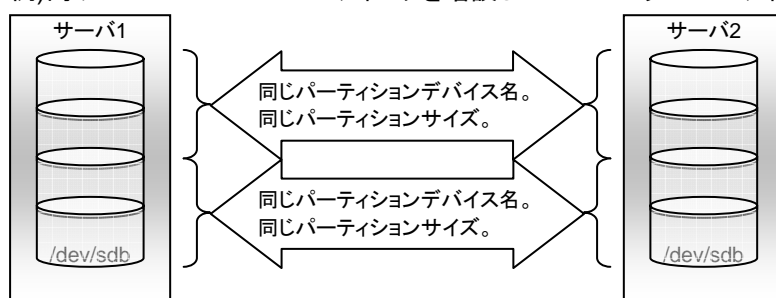
また 1 つのディスクに複数のミラーパーティションデバイスを割り当てて使用することができます。

例)両サーバに 2 つの SCSI ディスクを増設して 2 つのミラーディスクのペアにする場合。



- 1 つのディスク上にクラスタパーティションとデータパーティションをペアで確保してください。
- データパーティションを 1 つ目のディスク、クラスタパーティションを 2 つ目のディスクとするような使い方はできません。

例)両サーバに 1 つの SCSI ディスクを増設して 2 つのミラーパーティションにする場合



- ◆ ディスクに対してLinuxのmdやLVMによるストライプセット、ボリュームセット、ミラーリング、パリティ付きストライプセットの機能はサポートしていません。

依存するライブラリ

libxml2

OS インストール時に、libxml2 をインストールしてください。

依存するドライバ

softdog

- ◆ ユーザ空間モニタリソースの監視方法がsoftdogの場合、このドライバが必要です。
- ◆ ローダブルモジュール構成にしてください。スタティックドライバでは動作しません。

ミラードライバ

- ◆ ミラーパーティションのメジャー番号の218 を使用します。

他のデバイスドライバでは、メジャー番号の 218 を使用しないでください。

カーネルモードLANハートビートドライバ、キープアライブドライバ

- ◆ カーネルモードLANハートビートドライバは、メジャー番号 10、マイナ番号 240を使用します。
- ◆ キープアライブドライバは、メジャー番号 10、マイナ番号 241を使用します。

他のドライバが上記のメジャー及びマイナ番号を使用していないことを確認してください。

RAWモニタリソース用のパーティション確保

- ◆ RAWモニタリソースを設定する場合、監視専用のパーティションを用意してください。パーティションサイズは10MB確保してください。

OSインストール後、CLUSTERPROインストール前

OS のインストールが完了した後、OS やディスクの設定を行うときに留意頂いて頂きたいことです。

通信ポート番号

CLUSTERPRO では、デフォルトで以下のポート番号を使用します。このポート番号については「ミラードライブ間キープアライブ」以外は、Builder での変更が可能です。

下記ポート番号には、CLUSTERPRO 以外のプログラムからアクセスしないようにしてください。

サーバにファイアウォールの設定を行う場合には、下記のポート番号にアクセスできるようにしてください。

[サーバ・サーバ間]

From			To		備考
サーバ	自動割り当て ³	→	サーバ	29001/TCP	内部通信
サーバ	自動割り当て	→	サーバ	29002/TCP	データ転送
サーバ	自動割り当て	→	サーバ	29002/UDP	ハートビート
サーバ	自動割り当て	→	サーバ	29003/UDP	アラート同期
サーバ	自動割り当て	→	サーバ	29004/TCP	ミラーエージェント間通信
サーバ	自動割り当て	→	サーバ	29005/TCP	ミラードライブ間通信
サーバ	自動割り当て	→	サーバ	29006/UDP	ハートビート(カーネルモード)
サーバ	自動割り当て	→	サーバ	XXXX ⁴ /TCP	ミラーディスクリソースデータ同期
サーバ	自動割り当て	→	サーバ	icmp	ミラードライブ間キープアライブ

[サーバ・WebManager 間]

From			To		備考
Web マネージャ	自動割り当て	→	サーバ	29003/TCP	http 通信

[統合 WebManager を接続しているサーバ・管理対象のサーバ間]

From			To		備考
統合 WebManager を接続したサーバ	自動割り当て	→	サーバ	29003/TCP	http 通信

³ 自動割り当てでは、その時点で使用されていないポート番号が割り当てられます。

⁴ ミラーディスクリソースごとに使用するポート番号です。ミラーディスクリソース作成時に設定します。初期値として29051が設定されます。また、ミラーディスクリソースの追加ごとに1を加えた値が自動的に設定されます。変更する場合は、Builderの[ミラーディスクリソースプロパティ]-[詳細]タブで設定します。詳細については『リファレンスガイド』の第5章、「グループリソースの詳細」を参照してください。

時刻同期の設定

クラスタシステムでは、複数のサーバの時刻を定期的に同期する運用を推奨します。ntp などを使用してサーバの時刻を同期させてください。

NICデバイス名について

ifconfig コマンドの仕様により、CLUSTERPRO で動作可能な NIC デバイス名の文字列の長さに制限があります。また、フローティング IP リソースの数によって異なります。

NIC のデバイス名をデフォルト(eth0,eth1 など)から変更する場合には下記の範囲内の長さで設定してください。

bonding のデバイス名についても同様の制限があります。下記の "NIC デバイス名の文字列の長さ" の範囲内で bonding デバイスの名称を設定してください。

フローティングIPリソース の個数	NICデバイス名の文字列の長さ
0～10	7文字まで
11～100	6文字まで
100～	5文字まで

共有ディスクについて

- ◆ サーバの再インストール時等で共有ディスク上のデータを引き続き使用する場合は、パーティションの確保やファイルシステムの作成はしないでください。
- ◆ パーティションの確保やファイルシステムの作成をおこなうと共有ディスク上のデータは削除されます。
- ◆ 共有ディスク上のファイルシステムはCLUSTERPROが制御します。共有ディスクのファイルシステムをOSの/etc/fstabにエントリしないでください。
- ◆ 以下の手順で共有ディスクを設定します。

1. ディスクハートビート用パーティションの確保

- 共有ディスク上に CLUSTERPRO が独自に使用するパーティションを作成します。共有ディスクを使用するクラスタ内の 1 台のサーバから作成します。
- fdisk コマンドを使用してパーティションを確保します。パーティション ID は 83(Linux)で確保してください。
- 各ディスク(LUN)に 1 つディスクハートビートリソースで使用するパーティションを確保してください。
- ディスクハートビート用パーティションは 10MB(10x1024x1024 バイト)以上の大きさを確保してください(10MB ちょうどを指定しても、ディスクのジオメトリの違いにより実際には 10MB より大きなサイズが確保されますが、問題ありません)。

- 各 LUN にディスクハートビート専用パーティションを確保してください。ディスクの故障などでデバイス名がずれた場合にファイルシステムを破壊することがありますので、ディスクハートビートを使用しない LUN にもダミーのパーティションを確保してください。
- ディスクハートビート専用パーティションのパーティション番号が各 LUN で同じになるように確保してください。
- 複数の LUN を使用している場合でも、ディスクハートビートリソースはクラスタ内で 1 つまたは二つの使用を推奨します。ディスクハートビートリソースはハートビートインターバルごとにディスクへの read/write を行うためディスクへの負荷を考えて設定してください。

2. ディスクリソース用パーティションの確保

- 共有ディスク上にディスクリソースで使用するパーティションを作成します。共有ディスクを使用するクラスタ内の 1 台のサーバから作成します。
- fdisk コマンドを使用してパーティションを確保します。パーティション ID は 83(Linux)で確保してください。

3. ファイルシステムの作成

- 共有ディスク上のディスクリソース用パーティションにファイルシステムを構築します。
- 共有ディスクを使用するクラスタ内の 1 台のサーバから通常の Linux と同様に、mkfs コマンドなどでファイルシステムを構築してください。
- ディスクハートビート用パーティションにはファイルシステムの構築は必要ありません。
- 共有ディスクで使用するファイルシステムについて基本的に依存をしていますが、ファイルシステムの fsck の仕様により問題が発生することがあります。
- システムの対障害性の向上のために、ジャーナル機能を持つファイルシステムを使用することを推奨します。
- 現在 IA32、x86_64 で動作確認を完了しているファイルシステムは下記の通りです。

```
ext2
ext3
xfs
reiserfs
jfs
vxfs(対応 kernel に制限があります。
```

第 3 章 動作可能なディストリビューションと kernel を参照してください。)

- 現在 IA64 で動作確認を完了しているファイルシステムは下記の通りです。

```
ext2
ext3
xfs
```

4. マウントポイントの作成

- ディスクリソース用パーティションを mount するディレクトリを作成します。
- ディスクリソースを使用するクラスタ内のすべてのサーバで作成します。

ミラー用のディスクについて

- ◆ ミラーディスクリソース管理用パーティション(クラスタパーティション)とミラーディスクリソースで使用するパーティション(データパーティション)を設定します。
- ◆ ミラーディスク上のファイルシステムはCLUSTERPROが制御します。ミラーディスクのファイルシステムをOSの/etc/fstabにエントリしないでください。
- ◆ 以下の手順でミラーディスクを設定します。この手順は両方のサーバでおこないます。
 1. ミラーディスクの初期化(過去にCLUSTERPROのミラーディスクとして使用していたディスクを流用する場合のみ必要)
 - クラスタパーティションに以前のデータが残っているため初期化が必要です。
 - クラスタパーティションの初期化については「リファレンスガイド」を参照してください。
 2. クラスタパーティションの確保
 - ミラーディスク上に CLUSTERPRO が独自に使用するパーティションを作成します。
 - fdisk コマンドを使用してパーティションを確保します。
 - パーティション ID は 83(Linux)で確保してください。
 - 各ミラーディスクリソースに 1 つクラスタパーティションを確保してください。
 - クラスタパーティションは最低 10MB(10*1024*1024 バイト)の大きさを確保してください。
 - ディスクのジオメトリによっては 10MB 以上になる場合がありますが、問題ありません。
クラスタパーティションの詳細は「リファレンスガイド」を参照してください。
 3. データパーティションの確保
 - ミラーディスク上にミラーディスクリソースで使用するデータパーティションを作成します。
 - fdisk コマンドを使用してパーティションを確保します。
 - パーティション ID は 83(Linux)で確保してください。
 - データパーティションは 1GB 以上のサイズを確保してください。またパーティションサイズは 4096 バイトの倍数にしてください。ブロック数では 4 の倍数となります。
 - データパーティションの詳細は「リファレンスガイド」を参照してください。
 4. データパーティションのファイルシステムの作成
 - Builder でクラスタ構成情報作成時に、「初期 mkfs を行う」を設定する場合、CLUSTERPRO が自動でファイルシステムを構築します。
 - Builder でクラスタ構成情報作成時に、「初期 mkfs を行う」を選択しなければ CLUSTERPRO でファイルシステムの作成を行いません。
 - 「初期 mkfs を行う」の設定については「リファレンスガイド」を参照してください。

5. マウントポイントの作成

- ミラーディスクリソース用パーティションを mount するディレクトリを作成します。

OS起動時間の調整

電源が投入されてから、OS が起動するまでの時間が、下記の 2 つの時間より長くなるように調整してください。

- ◆ 共有ディスクを使用する場合に、ディスクの電源が投入されてから使用可能になるまでの時間
- ◆ ハートビートタイムアウト時間

OS ロードに lilo を使用している場合または GRUB を使用している場合の OS 起動時間の調整は、以下の手順になります。

lilo または GRUB 以外の OS ロードを使用している場合は、OS ロードの設定マニュアルを参照してください。

- ◆ liloを使用している場合

1. /etc/lilo.conf を編集します。

prompt オプションと timeout=<起動時間(単位は 1/10 秒)>オプションを指定します。または、prompt オプションを設定せず、delay=<起動時間(単位は 1/10 秒)>オプションを指定します。下記の例の場合にはアンダーラインの部分のみ変更してください。

```

---(例1, promptを出すケース, 起動時間30秒)---
boot=/dev/sda
map=/boot/map
install=/boot/boot.b
prompt
linear
timeout=300
image=/boot/vmlinuz
        label=linux
        root=/dev/sda1
        initrd=/boot/initrd.img
        read-only

---(例2, promptを出さないケース, 起動時間30秒)---
boot=/dev/sda
map=/boot/map
install=/boot/boot.b
#prompt
linear
delay=300
image=/boot/vmlinuz
        label=linux
        root=/dev/sda1
        initrd=/boot/initrd.img
        read-only

```

2. /sbin/lilo コマンドを実行して設定の変更を反映させます。

◆ GRUBを使用している場合

/boot/grub/menu.lst を編集します。

timeout <起動時間(単位は秒)> オプションを指定します。下記の例の場合にはアンダーラインの部分のみ変更してください。

```
---(例 起動時間30秒)---
default 0
timeout 30

title linux
    kernel (hd0,1)/boot/vmlinuz
    root=/dev/sda2    vga=785
    initrd (hd0,1)/boot/initrd

title floppy
    root (fd0)
    chainloader +1
```

ネットワークの確認

- ◆ インタコネクトやミラーディスクコネクトで使用するネットワークの確認をします。クラスタ内のすべてのサーバで確認します。
- ◆ ifconfigコマンドやpingコマンドを使用してネットワークの状態を確認してください。
 - public-LAN (他のマシンと通信を行う系)
 - インタコネクト専用 LAN(CLUSTERPRO のサーバ間接続に使用する系)
 - ミラーディスクコネクト LAN(インタコネクトと共用)
 - ホスト名
- ◆ クラスタで使用するフローティングIPリソースのIPアドレスは、OS側への設定は不要です。

ユーザ空間モニタリソース(監視方法ipmi)について

- ◆ 監視方法がipmiの場合、ipmiutilを使用します。
- ◆ CLUSTERPROにipmiutilは添付しておりません。ユーザ様ご自身で別途ipmiutil の rpm ファイルをインストールしてください。
- ◆ ipmiutilに関し以下の事項は弊社は対応いたしません。ユーザ様の判断、責任にてご使用ください。
 - ipmiutil 自体に関するお問い合わせ
 - ipmiutil の動作保証
 - ipmiutil の不具合対応、不具合が原因の障害
 - 各サーバの ipmiutil の対応状況のお問い合わせ
- ◆ ご使用予定のサーバ(ハードウェア)のipmiutil対応可否についてはユーザ様にて事前に確認ください。
- ◆ ハードウェアとしてIPMI規格に準拠している場合でも実際にはipmiutilが動作しない場合がありますので、ご注意ください。
- ◆ サーバベンダが提供するサーバ監視ソフトウェアを使用する場合には 監視方法にIPMIを選択しないでください。
これらのサーバ監視ソフトウェアとipmiutilは共にサーバ上のBMC(Baseboard Management Controller)を使用するため競合が発生して正しく監視が行うことができません。

ユーザ空間モニタリソース(監視方法softdog)について

- ◆ 監視方法にsoftdogを設定する場合、OS標準添付のheartbeatを動作しない設定にしてください。

CLUSTERPRO の情報作成時

CLUSTERPRO の構成情報の設計、作成前にシステムの構成に依存して確認、留意が必要な事項です。

グループリソースの非活性異常時の最終アクション

非活性異常検出時の最終動作に「何もしない」を選択すると、グループが非活性失敗のまま停止しません。

本番環境では「何もしない」は設定しないように注意してください。

VxVMが使用するRAWデバイスの確認

ボリューム RAW デバイスの実 RAW デバイスについて事前に調べておいてください。

1. CLUSTERPRO をインストールする前に、片サーバで活性しうる全てのディスクグループをインポートし、全てのボリュームを起動した状態にします。
2. 以下のコマンドを実行します。

```
# raw -qa
/dev/raw/raw2: bound to major 199, minor 2
/dev/raw/raw3: bound to major 199, minor 3
```

①②

例)ディスクグループ名、ボリューム名がそれぞれ以下の場合

- ディスクグループ名 dg1
- dg1 配下のボリューム名 vol1、vol2

3. 以下のコマンドを実行します。

```
# ls -l /dev/vx/dsk/dg1/
brw----- 1 root root 199, 2 5月 15 22:13 vol1
brw----- 1 root root 199, 3 5月 15 22:13 vol2
```

③

4. ②と③のメジャー/マイナ番号が等しいことを確認します。

これにより確認された RAW デバイス①は CLUSTERPRO のディスクハートビートリソース、RAW リソース、RAW モニタリソースでは絶対に設定しないでください。

ミラーディスクのファイルシステムの選択について

現在動作確認を完了しているファイルシステムは下記の通りです。

- ◆ ext2
- ◆ ext3
- ◆ xfs
- ◆ reiserfs
- ◆ jfs
- ◆ vxfs (対応kernelに関しては「第 3 章 CLUSTERPRO の動作環境」を参照してください。)

RAWモニタリソースについて

- ◆ RAWモニタリソースを設定する場合、既にmountしているパーティションまたはmountする可能性のあるパーティションの監視はできません。また、既にmountしているパーティションまたはmountする可能性のあるパーティションのwhole device(ディスク全体を示すデバイス)をデバイス名に設定して監視することもできません。
- ◆ 監視専用のパーティションを用意してRAWモニタリソースに設定してください。

遅延警告割合

遅延警告割合を 0 または、100 に設定すれば以下のようなことを行うことが可能です。

- ◆ 遅延警告割合に0を設定した場合
監視毎に遅延警告がアラート通報されます。
この機能を利用し、サーバが高負荷状態での監視リソースへのポーリング時間を算出し、監視リソースの監視タイムアウト時間を決定することができます。
- ◆ 遅延警告割合に100を設定した場合
遅延警告の通報を行いません。

テスト運用以外で、0%等の低い値を設定しないように注意してください。

ディスクモニタリソースの監視方法TURについて

- ◆ SCSIのTest Unit ReadyコマンドやSG_IOコマンドをサポートしていないディスク、ディスクインターフェイス(HBA)では使用できません。
ハードウェアがサポートしている場合でもドライバがサポートしていない場合があるのでドライバの仕様も合わせて確認してください。
- ◆ S-ATAインターフェイスのディスクの場合には、ディスクコントローラのタイプや使用するディストリビューションにより、OSにIDEインターフェイスのディスク(hd)として認識される場合とSCSIインターフェイスのディスク(sd)として認識される場合があります。
IDEインターフェイスとして認識される場合には、すべてのTUR方式は使用できません。
SCSIインターフェイスとして認識される場合には、TUR(legacy)が使用できます。
TUR(generic)は使用できません。
- ◆ Read方式に比べてOSやディスクへの負荷は小さくなります。
- ◆ Test Unit Readyでは、実際のメディアへのI/Oエラーは検出できない場合があります。

WebManagerの画面更新間隔について

- ◆ WebManagerタブの「画面データ更新インターバル」には、基本的に30秒より小さい値を設定しないでください。

LANハートビートの設定について

- ◆ LANハートビートリソースは最低一つ設定する必要があります。
- ◆ インタコネクト専用のLANをLANハートビートリソースとして登録し、さらにパブリックLANもLANハートビートリソースとして登録することを推奨します (LANハートビートリソースを二つ以上設定することを推奨します)。

カーネルモードLANハートビートの設定について

- ◆ カーネルモードLANハートビートが使用できるディストリビューション、カーネルの場合にはLANハートビートとカーネルモードLANハートビートの併用を推奨します。
- ◆ インタコネクト専用のLANをLANハートビートリソース、カーネルモードLANハートビートリソースとして登録し、さらにパブリックLANもLANハートビートリソース、カーネルモードLANハートビートリソースとして登録することを推奨します (LANハートビートリソースとカーネルモードLANハートビートリソースを二つ以上設定することを推奨します)。

COMハートビートの設定について

- ◆ ネットワークが断線した場合に両系で活性することを防ぐため、COMが使用できる環境であればCOMハートビートリソースを使用することを推奨します。

CLUSTERPRO運用後

クラスタとして運用を開始した後に発生する事象で留意して頂きたい事項です。

hotplugサービスについて

hotplug サービスがデバイスをサーチするときに以下のログが messages ファイルにエントリされます。

```
kernel: <liscal liscal_make_request> NMP0 I/O port is close,
mount(0), io(0).
kernel: Buffer I/O error on device NMP1, logical block 0
```

hotplug サービスが起動する時点でミラーリソースが起動していないため、この現象が発生します。この現象は異常ではありません。
hotplug を使用しない設定に変更して、coldplug で運用する場合には本現象は発生しません。

X-Window上のファイル操作ユーティリティについて

X-Window 上で動作する一部のファイル操作ユーティリティ(GUI でファイルやディレクトリのコピーや移動などの操作を行うもの)に以下の挙動をするものがあります。

- ◆ ブロックデバイスが使用可能であるかサーチする
- ◆ サーチの結果、マウントが可能なファイルシステムがあればマウントする

上記のような仕様のファイル操作ユーティリティは使用しないでください。
上記のような動作は CLUSTERPRO の動作に支障が発生する可能性があります。

ドライバロード時のメッセージについて

ミラードライバを load した際に、以下のメッセージがコンソール、syslog に表示されることがあります。この現象は異常ではありません。

```
kernel: liscal: no version for "struct_module" found: kernel
tainted.
kernel: liscal: module license 'unspecified' taints kernel.
```

clpka ドライバ、clpkhb ドライバを load した際に、以下のメッセージがコンソール、syslog に表示されることがあります。この現象は異常ではありません。

```
kernel: clpkhb: no version for "struct_module" found: kernel
tainted.
kernel: clpkhb: module license 'unspecified' taints kernel.

kernel: clpka: no version for "struct_module" found: kernel tainted.
kernel: clpka: module license 'unspecified' taints kernel.
```

ipmiのメッセージについて

ユーザ空間モニタリソースにIPMIを使用する場合、syslogに下記のkernelモジュール警告ログが多数出力されます。

```
modprobe: modprobe: Can't locate module char-major-10-173
```

このログ出力を回避したい場合は、/dev/ipmikcs を rename してください。

回復動作中の操作制限

モニタリソースの異常検出時の設定で回復対象にグループリソース(ディスクリソース、EXECリソース、...)を指定し、モニタリソースが異常を検出した場合の回復動作遷移中(再活性化 → フェイルオーバー → 最終動作)には、以下のコマンドまたは、WebManagerからのクラスタ及びグループへの制御は行わないでください。

- ◆ クラスタの停止 / サスペンド
- ◆ グループの開始 / 停止 / 移動

モニタリソース異常による回復動作遷移中に上記の制御を行うと、そのグループの他のグループリソースが停止しないことがあります。

また、モニタリソース異常状態であっても最終動作実行後であれば上記制御を行うことが可能です。

コマンド編に記載されていない実行形式ファイルやスクリプトファイルについて

インストールディレクトリ配下にコマンド編に記載されていない実行形式ファイルやスクリプトファイルがありますが、CLUSTERPRO 以外からは実行しないでください。

実行した場合の影響については、サポート対象外とします。

kernelページアロケートエラーのメッセージについて

TurboLinux 10 Server で Replicator を使用する場合に、syslogに以下のメッセージが出力されることがあります。ただし、物理メモリサイズや I/O 負荷に依存するため出力されない場合もあります。

```
kernel: [kernel モジュール名]: page allocation failure. order:X,  
mode:0xxx
```

このメッセージが出力される場合には、下記のカーネルパラメータを変更する必要があります。sysctl コマンド等を使用して OS 起動時にパラメータが変更されるように設定してください。

```
/proc/sys/vm/min_free_kbytes
```

min_free_kbytes に設定可能な最大値は、サーバに搭載されている物理メモリサイズによって異なります。下記の表を参照して設定してください。

物理メモリサイズ(Mbyte)	最大値
1024	1024
2048	1448
4096	2048
8192	2896
16384	4096

ログ収集時のメッセージ

ログ収集を実行した場合、コンソールに以下のメッセージが表示されることがありますが、異常ではありません。ログは正常に収集されています。

```
hd#: bad special flag: 0x03
ip_tables: (C) 2000-2002 Netfilter core team
```

(hd#にはサーバ上に存在する IDE のデバイス名が入ります)

クラスタシャットダウン・クラスタシャットダウンリブート

ミラーディスク使用時は、グループ活性処理中に clpstdn コマンドまたは WebManager からクラスタシャットダウン、クラスタシャットダウンリブートを実行しないでください。

グループ活性処理中はグループ非活性ができません。このため、ミラーディスクリソースが正常に非活性されていない状態で OS がシャットダウンされ、ミラーブレイクが発生することがあります。

特定サーバのシャットダウン、リブート

ミラーディスク使用時は、グループ活性処理中に clpdown コマンドまたは WebManager からサーバのシャットダウン、シャットダウンリブートコマンドを実行しないでください。

グループ活性処理中はグループ非活性ができません。このため、ミラーディスクリソースが正常に非活性されていない状態で OS がシャットダウンされ、ミラーブレイクが発生することがあります。

WebManagerについて

- ◆ WebManagerで表示される内容は必ずしも最新の状態を示しているわけではありません。最新の情報を取得したい場合、[リロード]ボタンを選択して最新の情報を取得してください。
- ◆ WebManagerが情報を取得中にサーバダウン等発生すると、情報の取得に失敗し、一部オブジェクトが正しく表示できない場合があります。次の自動更新まで待つか、[リロード]ボタンを選択して最新の情報を再取得してください。

- ◆ Linux上のブラウザを利用する場合、ウィンドウマネージャの組み合わせによっては、ダイアログが背後に回ってしまう場合があります。[ALT]+[TAB]キーなどでウィンドウを切り替えてください。
- ◆ CLUSTERPROのログ収集は複数のWebManagerから同時に実行することはできません。
- ◆ 接続先と通信できない状態で操作を行うと、制御が戻ってくるまでしばらく時間が必要な場合があります。
- ◆ マウスポインタが処理中を表す、腕時計や砂時計になっている状態で、ブラウザ外にカーソルを移動すると、処理中であってもカーソルが矢印の状態にもどってしまうことがあります。
- ◆ Proxyサーバを経由する場合は、WebManagerのポート番号を中継できるように、Proxyサーバの設定をしてください。
- ◆ CLUSTERPROのアップデートを行なった場合、ブラウザを終了してください。Javaのキャッシュをクリアしてブラウザを再起動してください。

Builder について

- ◆ 以下の製品とはクラスタ構成情報の互換性がありません。
 - CLUSTERPRO X 1.0 for Linux 以外の Linux 版のトレッキングツール
 - CLUSTERPRO for Windows Value Edition のトレッキングツール
- ◆ Webブラウザを終了すると(メニューの[終了]やウィンドウフレームの[X]ボタン等)、現在の編集内容が破棄されます。構成を変更した場合でも保存の確認ダイアログが表示されません。
編集内容の保存が必要な場合は、終了する前に、Builder のメニューバー[ファイル]-[情報ファイルの保存]を行ってください。
- ◆ Webブラウザをリロードすると(メニューの[最新の情報に更新]やツールバーの[現在のページを再読み込み]ボタン等)、現在の編集内容が破棄されます。構成を変更した場合でも保存の確認ダイアログが表示されません。
編集内容の保存が必要な場合は、リロードする前に、Builder のメニューバー[ファイル]-[情報ファイルの保存]を行ってください。

第 6 章 アップデート手順

本章では、CLUSTERPRO のアップデート手順について説明します。

本章で説明する項目は以下の通りです。

- CLUSTERPRO Ver3.x からのアップデート手順..... 70

CLUSTERPRO Ver3.x からのアップデート手順

クラスタ構成情報のバックアップ

クラスタ構成情報を FD にバックアップします。

クラスタ構成情報のバックアップは、root 権限を持つユーザで実行してください。以下の手順をマスタサーバで実行してください。

1. FD を装置にセットします。
2. FD が未フォーマット状態の場合には tar コマンドでできるように fdformat コマンドなどで format しておきます。
3. 以下のコマンドを実行します。

Linux 上で Builder を使用する場合

```
clpcfctrl --pull -l
```

Windows 上で Builder を使用する場合

```
clpcfctrl --pull -w
```

4. FD を装置から取り出して次の手順へ進みます。FD は下記の手順で使します。

クラスタ情報の変換

バックアップしたクラスタ情報を X 1.0 用のクラスタ情報に変換します。

Windows または Linux にインストールした X 1.0 用の Builder を使します。Builder のインストール手順については、インストールガイド&設計ガイドを参照してください。

1. FD を Builder を使用する PC またはサーバにセットします。
2. Builder を起動します。
3. Builder のメニューのファイル(F)→情報ファイルを開く(O)→クラスタ構成を変更(C)を実行します。
4. FD 上の clp.conf を指定して開きます。
5. Builder のメニューのファイル(F)→情報ファイルの保存(S)を実行します。上書き確認のダイアログに対してはい(Y)を選択します。
6. Builder を終了します。
7. FD を装置から取り出して次の手順へ進みます。FD は X 1.0 をインストールした後に使します。

3.xのアンインストール

アンインストールは、root 権限を持つユーザで実行してください。CLUSTERPRO Server は、以下の手順でアンインストールしてください。

1. `chkconfig --del name` を実行して以下の順序でサービスを無効にします。 *name* には以下のサービスを指定します。
 - clusterpro_alertsync
 - clusterpro_webmgr
 - clusterpro
 - clusterpro_md (LE の場合のみ)
 - clusterpro_trn
 - clusterpro_evt
2. サーバを再起動します。
3. `rpm -e clusterpro` を実行します。

X 1.0のインストール

CLUSTERPRO Server RPM は root ユーザでインストールしてください。次の手順に従って、サーバ RPM をすべてのサーバでインストールしてください。

1. インストール CD-ROM の媒体を mount します。
2. rpm コマンドを実行してパッケージファイルをインストールします。
アーキテクチャによりインストール用 RPM が異なります。

CD-ROM 内の `/Linux/1.0/jp/server` に移動して、

```
rpm -i clusterpro-<バージョン>.<アーキテクチャ>.rpm --nodeps
```

を実行します。

アーキテクチャには i686、x86_64、ia64、ppc64 があります。インストール先の環境に応じて選択してください。アーキテクチャは、`arch` コマンドなどで確認できます。

CLUSTERPRO は以下の場所にインストールされます。このディレクトリを変更するとアンインストールできなくなりますので注意してください。

インストールディレクトリ: `/opt/nec/clusterpro`

3. インストール終了後、インストール CD-ROM 媒体を umount します。
4. インストール CD-ROM 媒体を取り除いた後、サーバをリブートします。

X 1.0 のセットアップ手順に進みます。

付録

- 付録 A 用語集
- 付録 B 索引

付録 A 用語集

英数字

CLUSTERパーティション	ミラーディスクに設定するパーティション。ミラーディスクの管理に使用する。 関連(ディスクハートビート用パーティション)
----------------	--

あ

インタコネクト	クラスタ サーバ間の通信パス (関連) プライベート LAN、パブリック LAN
---------	---

か

仮想IPアドレス ⁵	遠隔地クラスタを構築する場合に使用するリソース (IPアドレス)
-----------------------	----------------------------------

管理クライアント	WebManager が起動されているマシン
----------	------------------------

起動属性	クラスタ起動時、自動的にフェイルオーバーグループを起動するか、手動で起動するかを決定するフェイルオーバー グループの属性 管理クライアントより設定が可能
------	---

共有ディスク	複数サーバよりアクセス可能なディスク
--------	--------------------

共有ディスク型クラスタ	共有ディスクを使用するクラスタシステム
-------------	---------------------

切替パーティション	複数のコンピュータに接続され、切り替えながら使用可能なディスクパーティション (関連)ディスクハートビート用パーティション
-----------	--

クラスタ システム	複数のコンピュータを LAN などをつないで、1 つのシステムのように振る舞わせるシステム形態
-----------	---

クラスタ シャットダウン	クラスタシステム全体 (クラスタを構成する全サーバ) をシャットダウンさせること
--------------	--

現用系	ある 1 つの業務セットについて、業務が動作しているサーバ (関連) 待機系
-----	---

⁵ 仮想IPアドレスはwindows版でのみ使用する概念になります。

さ

セカンダリ (サーバ)	通常運用時、フェイルオーバーグループがフェイルオーバーする先のサーバ (関連) プライマリ サーバ
--------------------	--

た

待機系	現用系ではない方のサーバ (関連) 現用系
ディスクハートビート用パーティション	共有ディスク型クラスターで、ハートビート通信に使用するためのパーティション
データパーティション	共有ディスクの切替パーティションのように使用することが可能なローカルディスク ミラーディスクに設定するデータ用のパーティション (関連) CLUSTER パーティション

な

ネットワークパーティション	全てのハートビートが途切れてしまうこと (関連) インタコネクト、ハートビート
ノード	クラスタシステムでは、クラスターを構成するサーバを指す。ネットワーク用語では、データを他の機器に経由することのできる、コンピュータやルータなどの機器を指す。

は

ハートビート	サーバの監視のために、サーバ間で定期的にお互いに通信を行うこと (関連) インタコネクト、ネットワークパーティション
パブリック LAN	サーバ / クライアント間通信パスのこと (関連) インタコネクト、プライベート LAN
フェイルオーバー	障害検出により待機系が、現用系上の業務アプリケーションを引き継ぐこと
フェイルバック	あるサーバで起動していた業務アプリケーションがフェイルオーバーにより他のサーバに引き継がれた後、業務アプリケーションを起動していたサーバに再び業務を戻すこと
フェイルオーバー グループ	業務を実行するのに必要なクラスタリソース、属性の集合

フェイルオーバー グループの移動	ユーザが意図的に業務アプリケーションを現用系から待機系に移動させること
フェイルオーバー ポリシー	フェイルオーバー可能なサーバリストとその中でのフェイルオーバー優先順位を持つ属性
プライベート LAN	クラスタを構成するサーバのみが接続された LAN (関連) インタコネクト、パブリック LAN
プライマリ (サーバ)	フェイルオーバーグループでの基準で主となるサーバ (関連) セカンダリ (サーバ)
フローティング IP アドレス	フェイルオーバーが発生したとき、クライアントのアプリケーションが接続先サーバの切り替えを意識することなく使用できる IP アドレス クラスタサーバが所属する LAN と同一のネットワーク アドレス内で、他に使用されていないホスト アドレスを割り当てる

ま

マスタサーバ	Builder の [クラスタのプロパティ]-[マスタサーバ] で先頭に表示されているサーバ
ミラーコネクト	データミラー型クラスタでデータのミラーリングを行うために使用する LAN。プライマリインタコネクトと兼用で設定することが可能。
ミラー ディスクシステム	共有ディスクを使用しないクラスタシステム サーバのローカルディスクをサーバ間でミラーリングする

付録 B 索引

B

Builder, 35, 42, 48, 67

C

CLUSTEREPRO, 17, 18
COMハートビート, 63

H

HA クラスタ, 4
hotplugサービス, 64

I

ipmiのメッセージ, 65

J

Java実行環境, 42, 43

K

kernel, 38
kernelページアロケートエラーのメッセージ, 65

L

LANハートビート, 63

N

NIC Link Up/Downモニタリソース, 50
NICデバイス名, 55

O

OS, 42, 43
OS起動時間, 58

R

RAWデバイス, 61
RAWモニタリソース, 53, 62

S

Single Point of Failure (SPOF), 3, 12

T

TUR, 62

W

WebManager, 35, 43, 48, 66
write性能, 50

あ

アップデート手順, 70
アプリケーションの引き継ぎ, 10
アンインストール, 71

い

依存するドライバ, 53
依存するライブラリ, 53
インストール, 71

か

カーネルモードLANハートビート, 63
カーネルモードLANハートビート、キープアライブドライバ, 53
画面更新間隔, 63
監視できる障害とできない障害, 21

き

業務監視, 20
共有ディスク, 55
共有ディスク要件, 49

く

クラスタオブジェクト, 28
クラスタ構成情報, 70
クラスタシステム, 3, 4
クラスタシャットダウン, 66
クラスタシャットダウンリブート, 66
クラスタリソースの引き継ぎ, 9
グループリソース, 29, 61

け

検出できる障害とできない障害, 21

さ

サーバ監視, 19
最終アクション, 61

し

資源, 23

時刻同期, 55
システム構成, 23
実行形式ファイル, 65
障害監視, 15, 19
障害検出, 3, 8

す

スクリプトファイル, 65
スペック, 36

せ

製品構成, 18

そ

ソフトウェア, 38
ソフトウェア構成, 18

ち

遅延警告割合, 62

つ

通信ポート番号, 54

て

ディスクインターフェイス, 36
ディスクサイズ, 41
ディスク容量, 42, 43
ディストリビューション, 38, 56
データの引き継ぎ, 9

と

動作OS, 48
特定サーバのシャットダウン, 66
特定サーバのシャットダウンリポート, 66
ドライバロード時のメッセージ, 64

な

内部監視, 20

ね

ネットワーク, 59
ネットワークインターフェイス, 37
ネットワークパーティション問題, 9

は

ハードウェア, 36
ハードウェア構成, 26, 27
ハートビートリソース, 29
バックアップ, 70

ふ

ファイルシステム, 51, 62
ファイル操作ユーティリティ, 64
フェイルオーバー, 11, 22
ブラウザ, 42, 43

み

ミラーディスク要件, 48
ミラードライバ, 53
ミラー用ディスク, 51, 57

め

メモリ容量, 41, 42, 43

も

モニタリソース, 30

ゆ

ユーザ空間モニタリソース, 60

り

リソース, 17, 29

ろ

ログ収集時のメッセージ, 66