

ExpressCluster[®] X 3.0 for Linux

Getting Started Guide

10/01/2010
First Edition



Revision History

Edition	Revised Date	Description
First	10/01/2010	New manual

© Copyright NEC Corporation 2010. All rights reserved.

Disclaimer

Information in this document is subject to change without notice. No part of this document may be reproduced or transmitted in any form by any means, electronic or mechanical, for any purpose, without the express written permission of NEC Corporation.

Trademark Information

ExpressCluster[®] X is a registered trademark of NEC Corporation.

FastSync[™] is a trademark of NEC Corporation.

Linux is a registered trademark or trademark of Linus Torvalds in the United States and other countries.

RPM is a trademark of Red Hat, Inc.

Intel, Pentium and Xeon are registered trademarks or trademarks of Intel Corporation.

Microsoft and Windows are registered trademarks of Microsoft Corporation in the United States and other countries.

Turbolinux is a registered trademark of Turbolinux, Inc.

VERITAS, VERITAS Logo and all other VERITAS product names and slogans are trademarks or registered trademarks of VERITAS Software Corporation.

Java is a trademark or registered trademark of Sun Microsystems, Inc. in the United States and other countries.

VMware is a registered trademark or trademark of VMware, Inc. in the United States and other countries.

Novell is a registered trademark of Novell, Inc. in the United States and Japan.

SUSE is a registered trademark of SUSE LINUX AG, a group company of U.S. Novell.

Citrix, Citrix XenServer, and Citrix Essentials are registered trademarks or trademarks of Citrix Systems, Inc. in the United State and other countries.

Other product names and slogans written in this manual are trademarks or registered trademarks of their respective companies.

Table of Contents

Preface	viii
Who Should Use This Guide.....	viii
How This Guide is Organized.....	viii
ExpressCluster X Documentation Set.....	ix
Conventions.....	x
Contacting NEC.....	xi
Section I Introducing ExpressCluster	13
Chapter 1 What is a cluster system?	15
Overview of the cluster system.....	16
High Availability (HA) cluster	16
Shared disk type.....	17
Data mirror type.....	19
Error detection mechanism.....	20
Problems with shared disk type.....	20
Network partition (split-brain-syndrome)	21
Taking over cluster resources	21
Taking over the data.....	21
Taking over the applications	22
Summary of failover	23
Eliminating single point of failure.....	24
Shared disk	24
Access path to the shared disk.....	25
LAN	26
Operation for availability.....	26
Failure monitoring.....	27
Chapter 2 Using ExpressCluster	29
What is ExpressCluster?	30
ExpressCluster modules.....	30
Software configuration of ExpressCluster	30
How an error is detected in ExpressCluster	32
What is server monitoring?.....	32
What is application monitoring?	33
What is internal monitoring?.....	33
Monitorable and non-monitorable errors.....	33
Detectable and non-detectable errors by server monitoring	33
Detectable and non-detectable errors by application monitoring	34
Network partition resolution.....	34
Failover mechanism.....	34
Failover resources	35
System configuration of the failover type cluster.....	36
Hardware configuration of the shared disk type cluster	39
Hardware configuration of the mirror disk type cluster	40
Hardware configuration of the hybrid disk type cluster	41
What is cluster object?	42
What is a resource?	43
Heartbeat resources	43
Network partition resolution resources	43
Group resources	43
Monitor resources	44
Getting started with ExpressCluster.....	47
Latest information.....	47
Designing a cluster system.....	47
Configuring a cluster system.....	47
Troubleshooting the problem.....	47
Section II Installing ExpressCluster	49

Chapter 3	Installation requirements for ExpressCluster	51
Hardware	52	
General server requirements	52	
Supported disk interfaces	52	
Supported network interfaces	53	
Software	54	
System requirements for ExpressCluster Server	54	
Supported distributions and kernel versions	54	
Applications supported by monitoring options	58	
Operation Environment of VM resources	60	
Required memory and disk size	61	
System requirements for the Builder	62	
Supported operating systems and browsers	62	
Java runtime environment	62	
Required memory and disk size	62	
Supported ExpressCluster versions	62	
System requirements for the WebManager	63	
Supported operating systems and browsers	63	
Java runtime environment	63	
Required memory and disk size	63	
System requirements for the Integrated WebManager	64	
Supported operating systems and browsers	64	
Java runtime environment	64	
Required memory size and disk size	64	
Chapter 4	Latest version information	65
Correspondence list of ExpressCluster and a manual	66	
Enhanced functions	67	
Corrected information	69	
Chapter 5	Notes and Restrictions	71
Designing a system configuration	72	
Function list and necessary license	72	
Supported operating systems for the Builder and WebManager	72	
Hardware requirements for mirror disks	72	
Hardware requirements for shared disks	73	
Hardware requirements for hybrid disks	75	
NIC link up/down monitor resource	76	
Write function of the mirror disk resource and hybrid disk resource	77	
Not outputting syslog to the mirror disk resource or the hybrid disk resource	77	
Notes when terminating the mirror disk resource or the hybrid disk resource	77	
Data consistency among multiple asynchronous mirror disks	78	
O_DIRECT for mirror or hybrid disk resources	78	
Installing operating system	78	
/opt/nec/clusterpro file system	78	
Mirror disks	79	
Hybrid disks	81	
Dependent library	81	
Dependent driver	82	
The major number of Mirror driver	82	
The major number of Kernel mode LAN heartbeat and keepalive drivers	82	
Partition for RAW monitoring of disk monitor resources	82	
SELinux settings	82	
Before installing ExpressCluster	83	
Communication port number	83	
Changing the range of automatic allocation for the communication port numbers	85	
Clock synchronization	85	
NIC device name	85	
Shared disk	85	
Mirror disk	85	
Hybrid disk	86	
Adjusting OS startup time	86	
Verifying the network settings	86	

Ipmitool and OpenIPMI	86
User mode monitor resource (monitoring method: softdog)	87
Log collection	87
nsupdate and nslookup	87
Notes when creating ExpressCluster configuration data	88
Environment variable	88
Force stop function, chassis identify lamp linkage	88
Server reset, server panic and power off	88
Final action for group resource deactivation error	89
Stack size of the application executed by EXEC resource	89
Verifying raw device for VxVM	89
Selecting mirror disk file system	90
Selecting hybrid disk file system	90
RAW monitoring of disk monitor resources	90
Delay warning rate	92
Disk monitor resource (monitoring method TUR)	92
WebManager reload interval	92
LAN heartbeat settings	92
Kernel mode LAN heartbeat resource settings	92
COM heartbeat resource settings	92
Double-byte character set that can be used in script comments	93
Failover exclusive attribute of virtual machine group	93
After start operating ExpressCluster	94
Error message in the load of the mirror driver in the udev environment	94
File operating utility on X-Window	94
Messages displayed when loading a driver	94
Messages written to syslog when multiple mirror disk resources or hybrid disk resources are used	95
IPMI message	96
Limitations during the recovery operation	96
Executable format file and script file not described in manuals	96
Message of kernel page allocation error	96
Cache swell by a massive I/O to mirror disk resource and hybrid disk resource	97
Messages when collecting logs	97
Cluster shutdown and reboot	98
Shutdown and reboot of individual server	98
Scripts for starting/stopping ExpressCluster services	98
Service startup time	99
Scripts in EXEC resources	99
Monitor resources that monitoring timing is “Active”	99
Notes on the WebManager	100
Notes on the Builder (Config mode of Cluster Manager)	100
Changing the partition size of mirror disks and hybrid disk resources	100
Chapter 6 Upgrading ExpressCluster	103
How to update from ExpressCluster X 2.0 or 2.1	104
How to upgrade from X2.0 or X2.1	104
Appendix A. Glossary	109
Appendix B. Index	111

Preface

Who Should Use This Guide

ExpressCluster Getting Started Guide is intended for first-time users of the ExpressCluster. The guide covers topics such as product overview of the ExpressCluster, how the cluster system is installed, and the summary of other available guides. In addition, latest system requirements and restrictions are described.

How This Guide is Organized

Section I **Introducing ExpressCluster**

Chapter 1 What is a cluster system?

Helps you to understand the overview of the cluster system and ExpressCluster.

Chapter 2 Using ExpressCluster

Provides instructions on how to use a cluster system and other related-information.

Section II **Installing ExpressCluster**

Chapter 3 Installation requirements for ExpressCluster

Provides the latest information that needs to be verified before starting to use ExpressCluster.

Chapter 4 Latest version information

Provides information on latest version of the ExpressCluster.

Chapter 5 Notes and Restrictions

Provides information on known problems and restrictions.

Chapter 6 Upgrading ExpressCluster

Provides instructions on how to update the ExpressCluster.

Appendix

Appendix A Glossary

Appendix B Index

ExpressCluster X Documentation Set

The ExpressCluster X manuals consist of the following four guides. The title and purpose of each guide is described below:

Getting Started Guide

This guide is intended for all users. The guide covers topics such as product overview, system requirements, and known problems.

Installation and Configuration Guide

This guide is intended for system engineers and administrators who want to build, operate, and maintain a cluster system. Instructions for designing, installing, and configuring a cluster system with ExpressCluster are covered in this guide.

Reference Guide

This guide is intended for system administrators. The guide covers topics such as how to operate ExpressCluster, function of each module, maintenance-related information, and troubleshooting. The guide is supplement to the *Installation and Configuration Guide*.

Integrated WebManager Administrator's Guide

This guide is intended for system administrators who manage cluster systems using ExpressCluster with Integrated WebManager, and also intended for system engineers who introduce Integrated WebManager. This guide describes detailed issues necessary for introducing Integrated WebManager in the actual procedures.

Conventions

In this guide, **Note**, **Important**, **Related Information** are used as follows:

Note:

Used when the information given is important, but not related to the data loss and damage to the system and machine.

Important:

Used when the information given is necessary to avoid the data loss and damage to the system and machine.

Related Information:

Used to describe the location of the information given at the reference destination.

The following conventions are used in this guide.

Convention	Usage	Example
Bold	Indicates graphical objects, such as fields, list boxes, menu selections, buttons, labels, icons, etc.	In User Name , type your name. On the File menu, click Open Database .
Angled bracket within the command line	Indicates that the value specified inside of the angled bracket can be omitted.	<code>clpstat -s[-h <i>host_name</i>]</code>
#	Prompt to indicate that a Linux user has logged in as root user.	<code># clpcl -s -a</code>
Monospace (courier)	Indicates path names, commands, system output (message, prompt, etc), directory, file names, functions and parameters.	<code>/Linux/3.0/en/server/</code>
Monospace bold (courier)	Indicates the value that a user actually enters from a command line.	Enter the following: <code># clpcl -s -a</code>
<i>Monospace italic</i> (courier)	Indicates that users should replace italicized part with values that they are actually working with.	<code>rpm -i expressclsbuilder-<version_number>- <release_number>.i686.rpm</code>

Contacting NEC

For the latest product information, visit our website below:

<http://www.nec.co.jp/pfsoft/clusterpro/clp/overseas.html>

Section I Introducing ExpressCluster

This section helps you to understand the overview of ExpressCluster and its system requirements.
This section covers:

- Chapter 1 What is a cluster system?
- Chapter 2 Using ExpressCluster

Chapter 1 What is a cluster system?

This chapter describes overview of the cluster system.

This chapter covers:

- Overview of the cluster system..... 16
- High Availability (HA) cluster 16
- Error detection mechanism 20
- Taking over cluster resources 21
- Eliminating single point of failure 24
- Operation for availability 26

Overview of the cluster system

A key to success in today's computerized world is to provide services without them stopping. A single machine down due to a failure or overload can stop entire services you provide with customers. This will not only result in enormous damage but also in loss of credibility you once enjoyed.

A cluster system is a solution to tackle such a disaster. Introducing a cluster system allows you to minimize the period during which operation of your system stops (down time) or to avoid system-down by load distribution.

As the word "cluster" represents, a cluster system is a system aiming to increase reliability and performance by clustering a group (or groups) of multiple computers. There are various types of cluster systems, which can be classified into the following three listed below. ExpressCluster is categorized as a high availability cluster.

High Availability (HA) Cluster

In this cluster configuration, one server operates as an active server. When the active server fails, a stand-by server takes over the operation. This cluster configuration aims for high-availability and allows data to be inherited as well. The high availability cluster is available in the shared disk type, data mirror type or remote cluster type.

Load Distribution Cluster

This is a cluster configuration where requests from clients are allocated to load-distribution hosts according to appropriate load distribution rules. This cluster configuration aims for high scalability. Generally, data cannot be taken over. The load distribution cluster is available in a load balance type or parallel database type.

High Performance Computing (HPC) Cluster

This is a cluster configuration where CPUs of all nodes are used to perform a single operation. This cluster configuration aims for high performance but does not provide general versatility.

Grid computing, which is one of the types of high performance computing that clusters a wider range of nodes and computing clusters, is a hot topic these days.

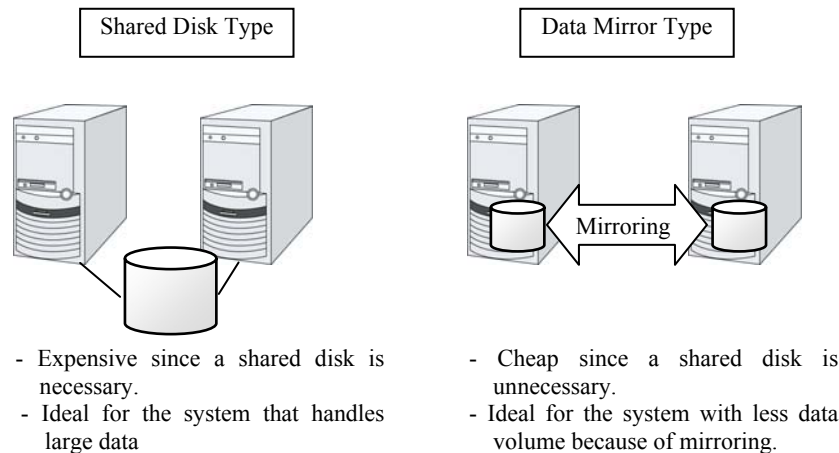
High Availability (HA) cluster

To enhance the availability of a system, it is generally considered that having redundancy for components of the system and eliminating a single point of failure is important. "Single point of failure" is a weakness of having a single computer component (hardware component) in the system. If the component fails, it will cause interruption of services. The high availability (HA) cluster is a cluster system that minimizes the time during which the system is stopped and increases operational availability by establishing redundancy with multiple servers.

The HA cluster is called for in mission-critical systems where downtime is fatal. The HA cluster can be divided into two types: shared disk type and data mirror type. The explanation for each type is provided below.

Shared disk type

Data must be inherited from one server to another in cluster systems. A cluster topology where data is stored in a shared disk with two or more servers using the data is called shared disk type.



- Expensive since a shared disk is necessary.
- Ideal for the system that handles large data

- Cheap since a shared disk is unnecessary.
- Ideal for the system with less data volume because of mirroring.

Figure 1-1: HA cluster configuration

If a failure occurs on a server where applications are running (active server), the cluster system detects the failure and applications are automatically started in a stand-by server to take over operations. This mechanism is called failover. Operations to be inherited in the cluster system consist of resources including disk, IP address and application.

In a non-clustered system, a client needs to access a different IP address if an application is restarted on a server other than the server where the application was originally running. In contrast, many cluster systems allocate a virtual IP address on an operational basis. A server where the operation is running, be it an active or a stand-by server, remains transparent to a client. The operation is continued as if it has been running on the same server.

File system consistency must be checked to inherit data. A check command (for example, `fsck` or `chkdsk` in Linux) is generally run to check file system consistency. However, the larger the file system is, the more time spent for checking. While checking is in process, operations are stopped. For this problem, journaling file system is introduced to reduce the time required for failover.

Logic of the data to be inherited must be checked for applications. For example, roll-back or roll-forward is necessary for databases. With these actions, a client can continue operation only by re-executing the SQL statement that has not been committed yet.

A server with the failure can return to the cluster system as a stand-by server if it is physically separated from the system, fixed, and then succeeds to connect the system. Such returning is acceptable in production environments where continuity of operations is important.

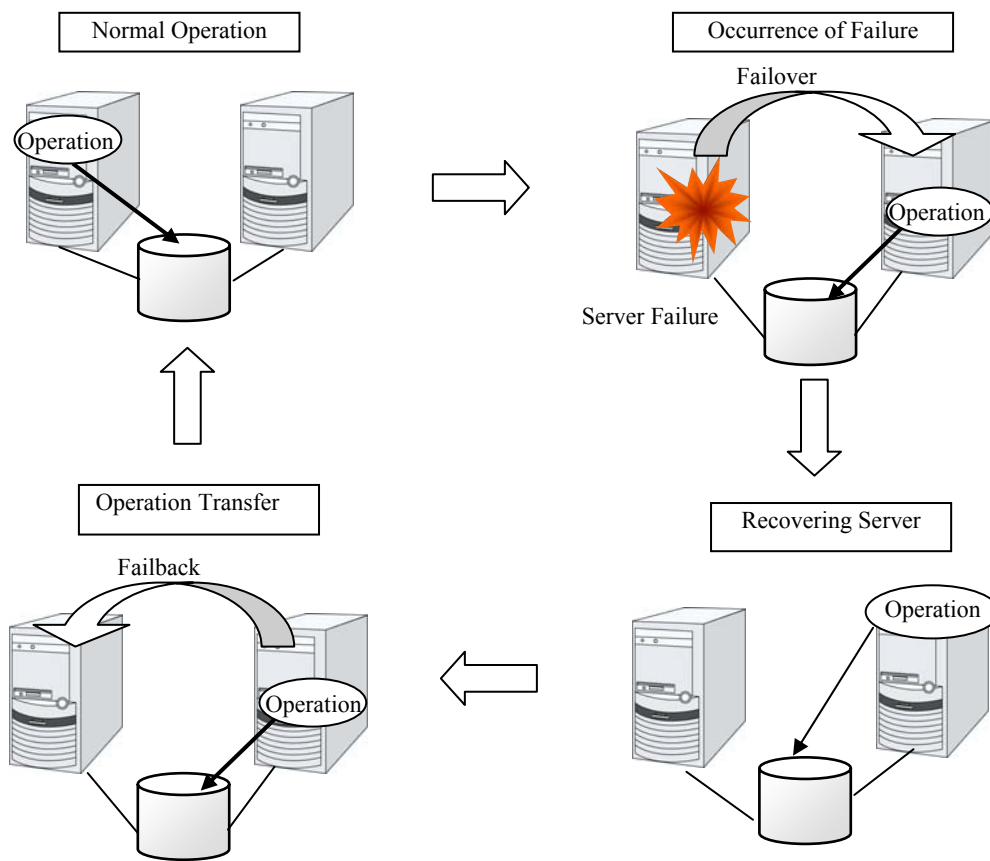


Figure 1-2: From occurrence of a failure to recovery

When the specification of the failover destination server does not meet the system requirements or overload occurs due to multi-directional stand-by, operations on the original server are preferred. In such a case, a failback takes place to resume operations on the original server.

A stand-by mode where there is one operation and no operation is active on the stand-by server, as shown in Figure 1-3, is referred to as uni-directional stand-by. A stand-by mode where there are two or more operations with each server of the cluster serving as both active and stand-by servers is referred to as multi-directional stand-by.

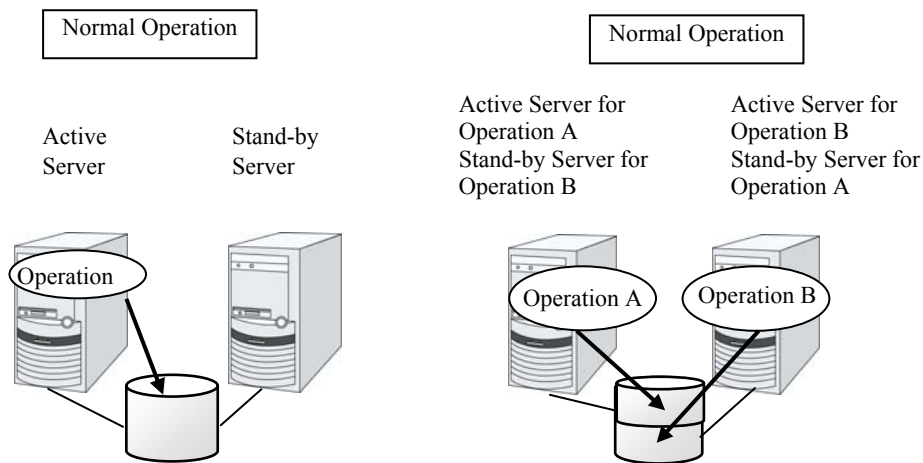


Figure 1-3: HA cluster topology

Data mirror type

The shared disk type cluster system is good for large-scale systems. However, creating a system with this type can be costly because shared disks are generally expensive. The data mirror type cluster system provides the same functions as the shared disk type with smaller cost through mirroring of server disks.

The data mirror type is not recommended for large-scale systems that handle a large volume of data since data needs to be mirrored between servers.

When a write request is made by an application, the data mirror engine not only writes data in the local disk but sends the write request to the stand-by server via the interconnect. Interconnect is a network connecting servers. It is used to monitor whether or not the server is activated in the cluster system. In addition to this purpose, interconnect is sometimes used to transfer data in the data mirror type cluster system. The data mirror engine on the stand-by server achieves data synchronization between stand-by and active servers by writing the data into the local disk of the stand-by server.

For read requests from an application, data is simply read from the disk on the active server.

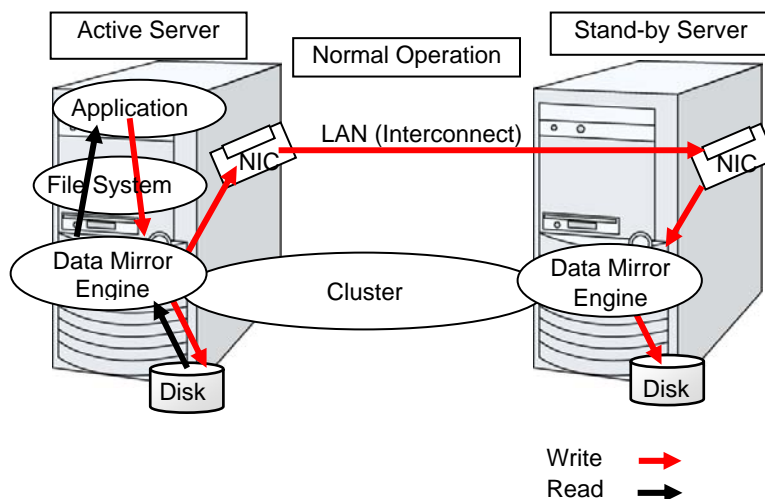


Figure 1-4: Data mirror mechanism

Snapshot backup is applied usage of data mirroring. Because the data mirror type cluster system has shared data in two locations, you can keep the disk of the stand-by server as snapshot backup without spending time for backup by simply separating the server from the cluster.

Failover mechanism and its problems

There are various cluster systems such as failover clusters, load distribution clusters, and high performance computing (HPC) clusters. The failover cluster is one of the high availability (HA) cluster systems that aim to increase operational availability through establishing server redundancy and passing operations being executed to another server when a failure occurs.

Error detection mechanism

Cluster software executes failover (for example, passing operations) when a failure that can impact continued operation is detected. The following section gives you a quick view of how the cluster software detects a failure.

Heartbeat and detection of server failures

Failures that must be detected in a cluster system are failures that can cause all servers in the cluster to stop. Server failures include hardware failures such as power supply and memory failures, and OS panic. To detect such failures, heartbeat is employed to monitor whether or not the server is active.

Some cluster software programs use heartbeat not only for checking whether or not the target is active through ping response, but for sending status information on the local server. Such cluster software programs begin failover if no heartbeat response is received in heartbeat transmission, determining no response as server failure. However, grace time should be given before determining failure, since a highly loaded server can cause delay of response. Allowing grace period results in a time lag between the moment when a failure occurred and the moment when the failure is detected by the cluster software.

Detection of resource failures

Factors causing stop of operations are not limited to stop of all servers in the cluster. Failure in disks used by applications, NIC failure, and failure in applications themselves are also factors that can cause the stop of operations. These resource failures need to be detected as well to execute failover for improved availability.

Accessing a target resource is a way employed to detect resource failures if the target is a physical device. For monitoring applications, trying to service ports within the range not impacting operation is a way of detecting an error in addition to monitoring whether or not application processes are activated.

Problems with shared disk type

In a failover cluster system of the shared disk type, multiple servers physically share the disk device. Typically, a file system enjoys I/O performance greater than the physical disk I/O performance by keeping data caches in a server.

What if a file system is accessed by multiple servers simultaneously?

Since a general file system assumes no server other than the local updates data on the disk, inconsistency between caches and the data on the disk arises. Ultimately the data will be corrupted. The failover cluster system locks the disk device to prevent multiple servers from mounting a file system, simultaneously caused by a network partition.

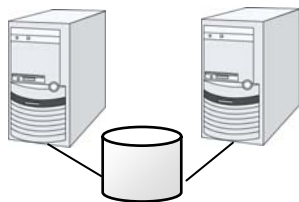


Figure 1-5: Cluster configuration with a shared disk

Network partition (split-brain-syndrome)

When all interconnects between servers are disconnected, failover takes place because the servers assume other server(s) are down. To monitor whether the server is activated, a heartbeat communication is used. As a result, multiple servers mount a file system simultaneously causing data corruption. This explains the importance of appropriate failover behavior in a cluster system at the time of failure occurrence.

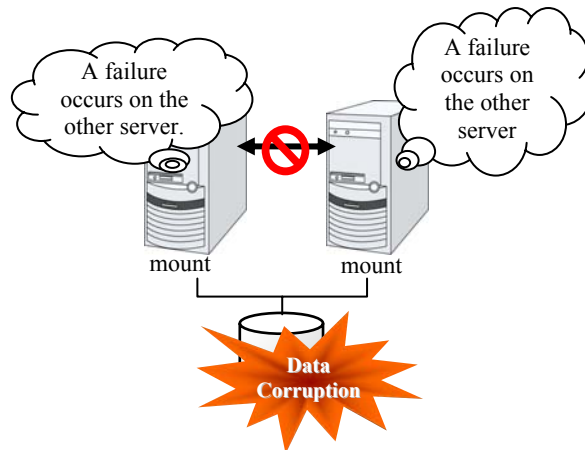


Figure 1-6: Network partition problem

The problem explained in the section above is referred to as “network partition” or “split-brain syndrome.” The failover cluster system is equipped with various mechanisms to ensure shared disk lock at the time when all interconnects are disconnected.

Taking over cluster resources

As mentioned earlier, resources to be managed by a cluster include disks, IP addresses, and applications. The functions used in the failover cluster system to inherit these resources are described below.

Taking over the data

Data to be passed from a server to another in a cluster system is stored in a partition on the shared disk. This means data is re-mounting the file system of files that the application uses on a healthy server. What the cluster software should do is simply mount the file system because the shared disk is physically connected to a server that inherits data.

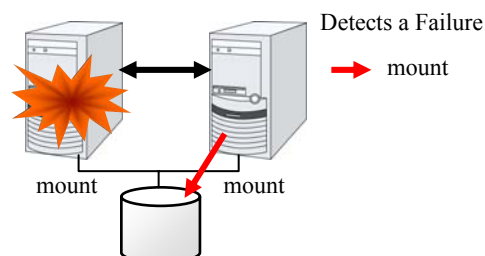


Figure 1-7: Taking over data

The figure 1-7 may look simple, but consider the following issues in designing and creating a cluster system.

One issue to consider is recovery time for a file system. A file system to be inherited may have been used by another server or being updated just before the failure occurred and requires a file system consistency check. When the file system is large, the time spent for checking consistency will be enormous. It may take a few hours to complete the check and the time is wholly added to the time for failover (time to take over operation), and this will reduce system availability.

Another issue you should consider is writing assurance. When an application writes important data into a file, it tries to ensure the data to be written into a disk by using a function such as synchronized writing. The data that the application assumes to have been written is expected to be inherited after failover. For example, a mail server reports the completion of mail receiving to other mail servers or clients after it has securely written mails it received in a spool. This will allow the spooled mail to be distributed again after the server is restarted. Likewise, a cluster system should ensure mails written into spool by a server to become readable by another server.

Taking over the applications

The last to come in inheritance of operation by cluster software is inheritance of applications. Unlike fault tolerant computers (FTC), no process status such as contents of memory is inherited in typical failover cluster systems. The applications running on a failed server are inherited by rerunning them on a healthy server.

For example, when instances of a database management system (DBMS) are inherited, the database is automatically recovered (roll-forward/roll-back) by startup of the instances. The time needed for this database recovery is typically a few minutes though it can be controlled by configuring the interval of DBMS checkpoint to a certain extent.

Many applications can restart operations by re-execution. Some applications, however, require going through procedures for recovery if a failure occurs. For these applications, cluster software allows to start up scripts instead of applications so that recovery process can be written. In a script, the recovery process, including cleanup of files half updated, is written as necessary according to factors for executing the script and information on the execution server.

Summary of failover

To summarize the behavior of cluster software:

- ◆ Detects a failure (heartbeat/resource monitoring)
- ◆ Resolves a network partition (NP resolution)
- ◆ Switches cluster resources
 - Pass data
 - Pass IP address
 - Application Taking over

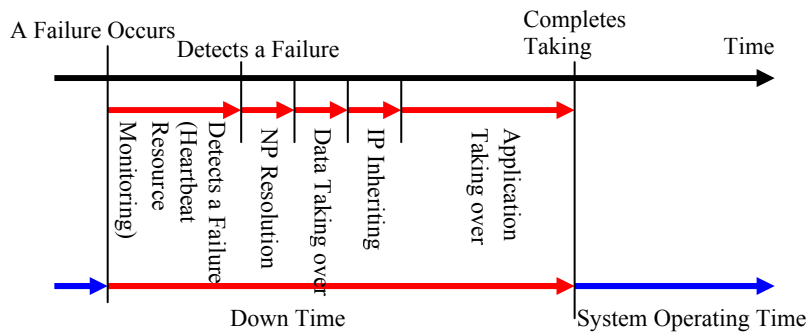


Figure 1-8: Failover time chart

Cluster software is required to complete each task quickly and reliably (see Figure 1-8.) Cluster software achieves high availability with due consideration on what has been described so far.

Eliminating single point of failure

Having a clear picture of the availability level required or aimed is important in building a high availability system. This means when you design a system, you need to study cost effectiveness of countermeasures, such as establishing a redundant configuration to continue operations and recovering operations within a short period of time, against various failures that can disturb system operations.

Single point of failure (SPOF), as described previously, is a component where failure can lead to stop of the system. In a cluster system, you can eliminate the system's SPOF by establishing server redundancy. However, components shared among servers, such as shared disk may become a SPOF. The key in designing a high availability system is to duplicate or eliminate this shared component.

A cluster system can improve availability but failover will take a few minutes for switching systems. That means time for failover is a factor that reduces availability. Solutions for the following three, which are likely to become SPOF, will be discussed hereafter although technical issues that improve availability of a single server such as ECC memory and redundant power supply are important.

- ◆ Shared disk
- ◆ Access path to the shared disk
- ◆ LAN

Shared disk

Typically a shared disk uses a disk array for RAID. Because of this, the bare drive of the disk does not become SPOF. The problem is the RAID controller is incorporated. Shared disks commonly used in many cluster systems allow controller redundancy.

In general, access paths to the shared disk must be duplicated to benefit from redundant RAID controller. There are still things to be done to use redundant access paths in Linux (described later in this chapter). If the shared disk has configuration to access the same logical disk unit (LUN) from duplicated multiple controllers simultaneously, and each controller is connected to one server, you can achieve high availability by failover between nodes when an error occurs in one of the controllers.

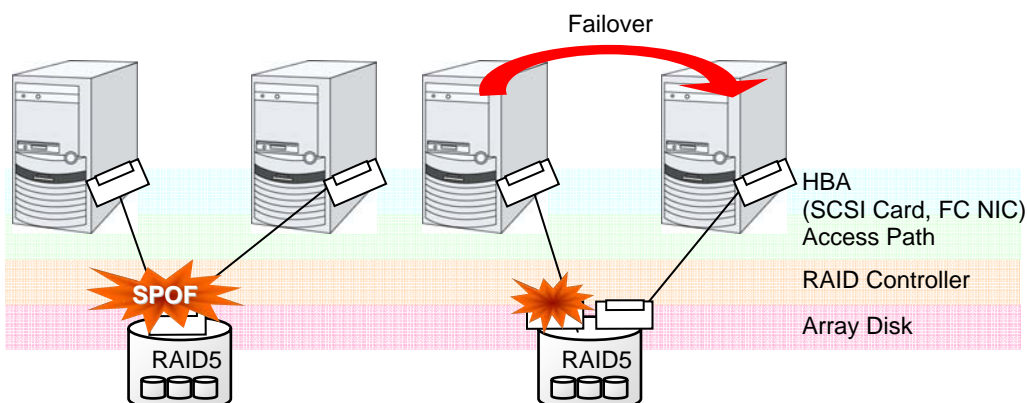


Figure 1-9: Example of the shared disk RAID controller and access paths being SPOF (left) and an access path connected to a RAID controller

With a failover cluster system of data mirror type, where no shared disk is used, you can create an ideal system having no SPOF because all data is mirrored to the disk in the other server. However you should consider the following issues:

- ◆ Disk I/O performance in mirroring data over the network (especially writing performance)
- ◆ System performance during mirror resynchronization in recovery from server failure (mirror copy is done in the background)
- ◆ Time for mirror resynchronization (clustering cannot be done until mirror resynchronization is completed)

In a system with frequent data viewing and a relatively small volume of data, choosing the data mirror type for clustering is a key to increase availability.

Access path to the shared disk

In a typical configuration of the shared disk type cluster system, the access path to the shared disk is shared among servers in the cluster. To take SCSI as an example, two servers and a shared disk are connected to a single SCSI bus. A failure in the access path to the shared disk can stop the entire system.

What you can do for this is to have a redundant configuration by providing multiple access paths to the shared disk and make them look as one path for applications. The device driver allowing such is called a path failover driver. Path failover drivers are often developed and released by shared disk vendors. Path failover drivers in Linux are still under development. For the time being, as discussed earlier, offering access paths to the shared disk by connecting a server on an array controller on the shared disk basis is the way to ensure availability in Linux cluster systems.

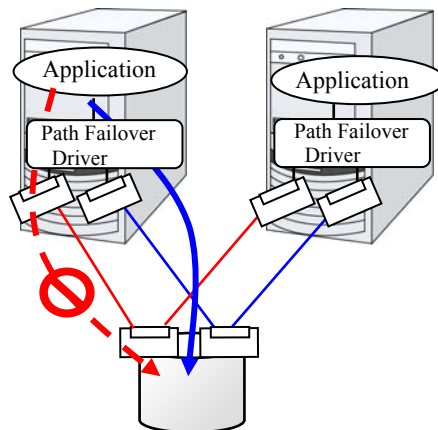


Figure 1-10: Path failover driver

LAN

In any systems that run services on a network, a LAN failure is a major factor that disturbs operations of the system. If appropriate settings are made, availability of cluster system can be increased through failover between nodes at NIC failures. However, a failure in a network device that resides outside the cluster system disturbs operation of the system.

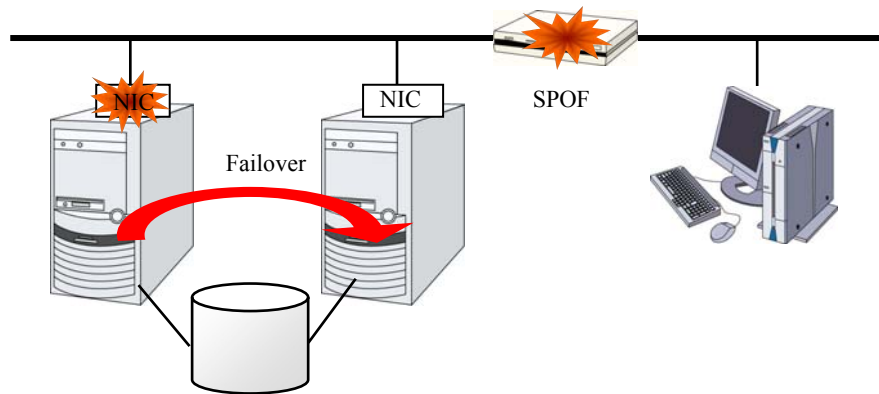


Figure 1-11: Example of router becoming SPOF

LAN redundancy is a solution to tackle device failure outside the cluster system and to improve availability. You can apply ways used for a single server to increase LAN availability. For example, choose a primitive way to have a spare network device with its power off, and manually replace a failed device with this spare device. Choose to have a multiplex network path through a redundant configuration of high-performance network devices, and switch paths automatically. Another option is to use a driver that supports NIC redundant configuration such as Intel's ANS driver.

Load balancing appliances and firewall appliances are also network devices that are likely to become SPOF. Typically they allow failover configurations through standard or optional software. Having redundant configuration for these devices should be regarded as requisite since they play important roles in the entire system.

Operation for availability

Evaluation before starting operation

Given many of factors causing system troubles are said to be the product of incorrect settings or poor maintenance, evaluation before actual operation is important to realize a high availability system and its stabilized operation. Exercising the following for actual operation of the system is a key in improving availability:

- ◆ Clarify and list failures, study actions to be taken against them, and verify effectiveness of the actions by creating dummy failures.
- ◆ Conduct an evaluation according to the cluster life cycle and verify performance (such as at degenerated mode)
- ◆ Arrange a guide for system operation and troubleshooting based on the evaluation mentioned above.

Having a simple design for a cluster system contributes to simplifying verification and improvement of system availability.

Failure monitoring

Despite the above efforts, failures still occur. If you use the system for long time, you cannot escape from failures: hardware suffers from aging deterioration and software produces failures and errors through memory leaks or operation beyond the originally intended capacity. Improving availability of hardware and software is important yet monitoring for failure and troubleshooting problems is more important. For example, in a cluster system, you can continue running the system by spending a few minutes for switching even if a server fails. However, if you leave the failed server as it is, the system no longer has redundancy and the cluster system becomes meaningless should the next failure occur.

If a failure occurs, the system administrator must immediately take actions such as removing a newly emerged SPOF to prevent another failure. Functions for remote maintenance and reporting failures are very important in supporting services for system administration. Linux is known for providing good remote maintenance functions. Mechanism for reporting failures are coming in place. To achieve high availability with a cluster system, you should:

- ◆ Remove or have complete control on single point of failure.
- ◆ Have a simple design that has tolerance and resistance for failures, and be equipped with a guide for operation and troubleshooting.
- ◆ Detect a failure quickly and take appropriate action against it.

Chapter 2 Using ExpressCluster

This chapter explains the components of ExpressCluster, how to design a cluster system, and how to use ExpressCluster.

This chapter covers:

- What is ExpressCluster?..... 30
- ExpressCluster modules 30
- Software configuration of ExpressCluster 30
- Network partition resolution..... 34
- Failover mechanism 34
- What is a resource? 43
- Getting started with ExpressCluster 47

What is ExpressCluster?

ExpressCluster is software that enhances availability and expandability of systems by a redundant (clustered) system configuration. The application services running on the active server are automatically inherited to a standby server when an error occurs in the active server.

ExpressCluster modules

ExpressCluster consists of following three modules:

ExpressCluster Server

A core component of ExpressCluster. Includes all high availability function of the server. The server function of the WebManager is also included.

ExpressCluster X WebManager (WebManager)

A tool to manage ExpressCluster operations. Uses a Web browser as a user interface. The WebManager is installed in ExpressCluster Server, but it is distinguished from the ExpressCluster Server because the WebManager is operated from the Web browser on the management PC.

ExpressCluster X Builder (Builder)

A tool for editing the cluster configuration data. The Builder also uses Web browser as a user interface. The following two versions of Builder are provided: the offline version, which is installed on your terminal as software independent of ExpressCluster Server, and the online version, which is opened by clicking the **setup mode** icon on the WebManager screen toolbar or **Setup Mode** on the **View** menu. The Builder needs to be installed separately from the ExpressCluster Server on the machine where you use the Builder.

Software configuration of ExpressCluster

The software configuration of ExpressCluster should look similar to the figure below. Install the ExpressCluster Server (software) on a Linux server, and the Builder on a management PC or a server. Because the main functions of WebManager and Builder are included in ExpressCluster Server, it is not necessary to separately install them. However, to use the Builder in an environment where ExpressCluster Server is not accessible, the offline version of Builder must be installed on the PC. The WebManager or Builder can be used through the Web browser on the management PC or on each server in the cluster.

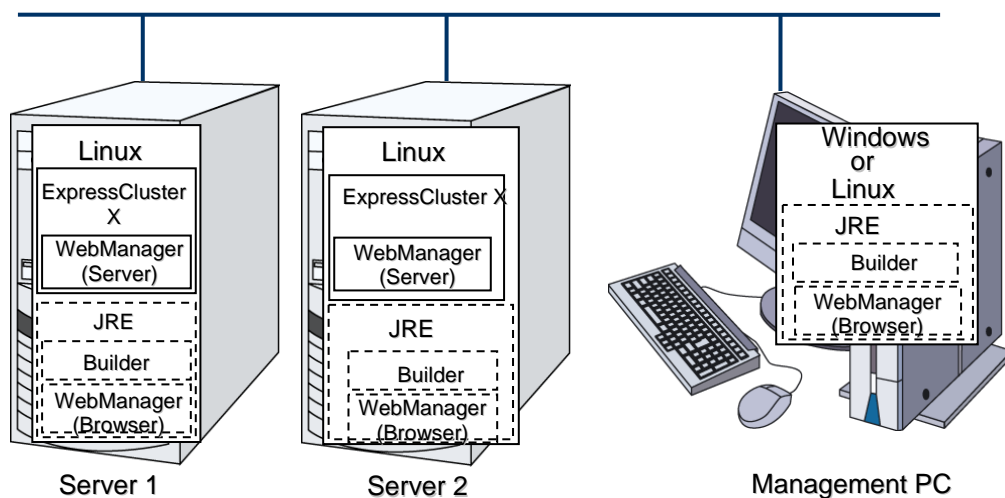


Figure 2-1 Software configuration of ExpressCluster

How an error is detected in ExpressCluster

There are three kinds of monitoring in ExpressCluster: (1) server monitoring, (2) application monitoring, and (3) internal monitoring. These monitoring functions let you detect an error quickly and reliably. The details of the monitoring functions are described below.

What is server monitoring?

Server monitoring is the most basic function of the failover-type cluster system. It monitors if a server that constitutes a cluster is properly working.

ExpressCluster regularly checks whether other servers are properly working in the cluster system. This way of verification is called “heartbeat communication.” The heartbeat communication uses the following communication paths:

Interconnect-dedicated LAN

Uses an Ethernet NIC in communication path dedicated to the failover-type cluster system. This is used to exchange information between the servers as well as to perform heartbeat communication.

Public LAN

Uses a communication path used for communication with client machine as an alternative interconnect. Any Ethernet NIC can be used as long as TCP/IP can be used. This is also used to exchange information between the servers and to perform heartbeat communication.

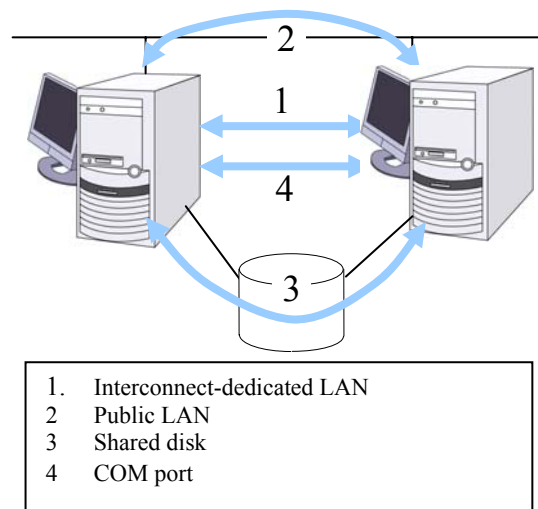
Shared disk

Creates an ExpressCluster-dedicated partition (ExpressCluster partition) on the disk that is connected to all servers that constitute the failover-type cluster system, and performs heartbeat communication on the ExpressCluster partition.

COM port

Performs heartbeat communication between the servers that constitute the failover-type cluster system through a COM port, and checks whether other servers are working properly.

Having these communication paths dramatically improves the reliability of the communication between the servers, and prevents the occurrence of network partition.



Note:

Network partition (also known as “split-brain syndrome”) refers to a condition when a network gets split by having a problem in all communication paths of the servers in a cluster. In a cluster system that is not capable of handling a network partition, a problem occurred in a communication path and a server cannot be distinguished. As a result, multiple servers may access the same resource and cause the data in a cluster system to be corrupted.

What is application monitoring?

Application monitoring is a function that monitors applications and factors that cause a situation where an application cannot run.

Activation status of application monitoring

An error can be detected by starting up an application from an exec resource in ExpressCluster and regularly checking whether a process is active or not by using the pid monitor resource. It is effective when the factor for application to stop is due to error termination of an application.

Note:

An error in resident process cannot be detected in an application started up by ExpressCluster. When the monitoring target application starts and stops a resident process, an internal application error (such as application stalling, result error) cannot be detected.

Resource monitoring

An error can be detected by monitoring the cluster resources (such as disk partition and IP address) and public LAN using the monitor resources of the ExpressCluster. It is effective when the factor for application to stop is due to an error of a resource which is necessary for an application to operate.

What is internal monitoring?

Internal monitoring refers to an inter-monitoring of modules within ExpressCluster. It monitors whether each monitoring function of ExpressCluster is properly working. Activation status of ExpressCluster process monitoring is performed within ExpressCluster.

Monitorable and non-monitorable errors

There are monitorable and non-monitorable errors in ExpressCluster. It is important to know what can or cannot be monitored when building and operating a cluster system.

Detectable and non-detectable errors by server monitoring

Monitoring condition: A heartbeat from a server with an error is stopped

Example of errors that can be monitored:

- ◆ Hardware failure (of which OS cannot continue operating)
- ◆ System panic

Example of error that cannot be monitored:

- ◆ Partial failure on OS (for example, only a mouse or keyboard does not function)

Detectable and non-detectable errors by application monitoring

Monitoring conditions: Termination of applications with errors, continuous resource errors, and disconnection of a path to the network devices.

Example of errors that can be monitored:

- ◆ Abnormal termination of an application
- ◆ Failure to access the shared disk (such as HBA¹ failure)
- ◆ Public LAN NIC problem

Example of errors that cannot be monitored:

- ◆ Application stalling and resulting in error. ExpressCluster cannot monitor application stalling and error results. However, it is possible to perform failover by creating a program that monitors applications and terminates itself when an error is detected, starting the program using the exec resource, and monitoring application using the PID monitor resource.

Network partition resolution

When the stop of a heartbeat is detected from a server, ExpressCluster determines whether it is an error in a server or a network partition. If it is judged as a server failure, failover (activate resources and start applications on a healthy server) is performed. If it is judged as network partition, protecting data is given priority over Taking over operations, so processing such as emergency shutdown is performed.

The following is the network partition resolution method:

- ◆ ping method

Related Information:

For the details on the network partition resolution method, see Chapter 7, “Details on network partition resolution resources” in Section II of the Reference Guide.

Failover mechanism

When an error is detected, ExpressCluster determines whether an error detected before failing over is an error in a server or a network partition. Then a failover is performed by activating various resources and starting up applications on a properly working server.

The group of resources which fail over at the same time is called a “failover group.” From a user’s point of view, a failover group appears as a virtual computer.

Note:

In a cluster system, a failover is performed by restarting the application from a properly working node. Therefore, what is saved in an application memory cannot be failed over.

From occurrence of error to completion of failover takes a few minutes. See the figure 2-2 below:

¹ HBA is an abbreviation for host bus adapter. This adapter is not for the shared disk, but for the server.

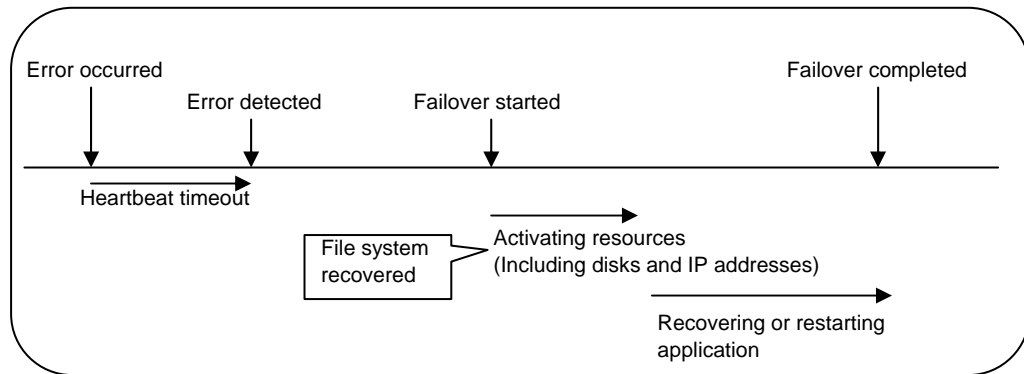


Figure 2-2 Failover time chart

Heartbeat timeout

- ◆ The time for a standby server to detect an error after that error occurred on the active server.
- ◆ The setting values of the cluster properties should be adjusted depending on the application load. (The default value is 90 seconds.)

Activating various resources

- ◆ The time to activate the resources necessary for operating an application.
- ◆ The resources can be activated in a few seconds in ordinary settings, but the required time changes depending on the type and the number of resources registered to the failover group. For more information, refer to the *Installation and Configuration Guide*.

Start script execution time

- ◆ The data recovery time for a roll-back or roll-forward of the database and the startup time of the application to be used in operation.
- ◆ The time for roll-back or roll-forward can be predicted by adjusting the check point interval. For more information, refer to the document that comes with each software product.

Failover resources

ExpressCluster can fail over the following resources:

Switchable partition

- ◆ Resources such as disk resource, mirror disk resource and hybrid disk resource.
- ◆ A disk partition to store the data that the application takes over.

Floating IP Address

- ◆ By connecting an application using the floating IP address, a client does not have to be conscious about switching the servers due to failover processing.
- ◆ It is achieved by dynamic IP address allocation to the public LAN adapter and sending ARP packet. Connection by floating IP address is possible from most of the network devices.

Script (exec resource)

- ◆ In ExpressCluster, applications are started up from the scripts.
- ◆ The file failed over on the shared disk may not be complete as data even if it is properly working as a file system. Write the recovery processing specific to an application at the time of failover in addition to the startup of an application in the scripts.

Note:

In a cluster system, failover is performed by restarting the application from a properly working node. Therefore, what is saved in an application memory cannot be failed over.

System configuration of the failover type cluster

In a failover-type cluster, a disk array device is shared between the servers in a cluster. When an error occurs on a server, the standby server takes over the applications using the data on the shared disk.

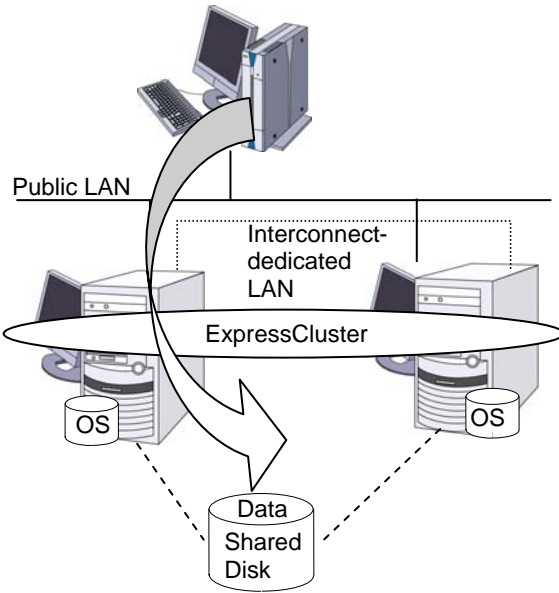


Figure 2-3 System configuration

A failover-type cluster can be divided into the following categories depending on the cluster topologies:

Uni-Directional Standby Cluster System

In the uni-directional standby cluster system, the active server runs applications while the other server, the standby server, does not. This is the simplest cluster topology and you can build a high-availability system without performance degradation after failing over.

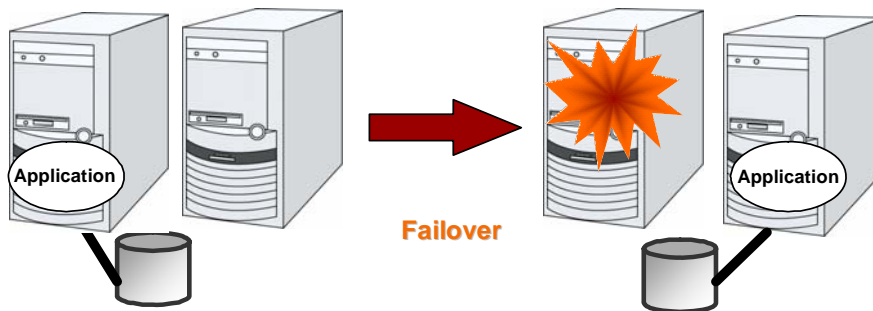
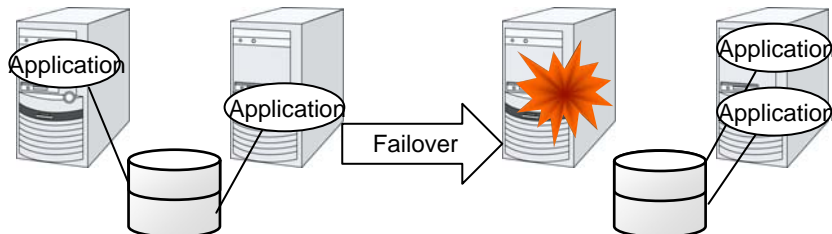


Figure 2-4 Uni-directional standby cluster system

Same Application Multi Directional Standby Cluster System

In the same application multi-directional standby cluster system, the same applications are activated on multiple servers. These servers also operate as standby servers. The applications must support multi-directional standby operation. When the application load data can be split into multiple data, depending on the data to be accessed, you can build a load distribution system per data partitioning basis by changing the client's connecting server.

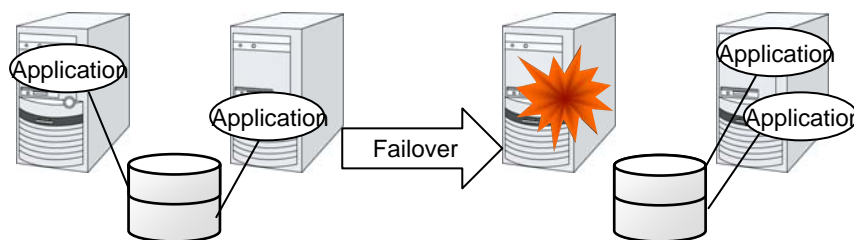


- The applications in the diagram are the same application.
- Multiple application instances are run on a single server after failover.

Figure 2-5 Same application multi directional standby cluster system

Different Application – Multi Directional Standby Cluster System

In the different application multi-directional standby cluster system, different applications are activated on multiple servers and these servers also operate as standby servers. The applications do not have to support multi-directional standby operation. A load distribution system can be built per application unit basis.



- Operation 1 and operation 2 use different applications.

Figure 2-6 Different application multi directional standby cluster system

Node to Node Configuration

The configuration can be expanded with more nodes by applying the configurations introduced thus far. In a node to node configuration described below, three different applications are run on three servers and one standby server takes over the application if any problem occurs. In a uni-directional standby cluster system, one of the two servers functions as a standby server. However, in a node to node configuration, only one of the four server functions as a standby server and performance deterioration is not anticipated if an error occurs only on one server.

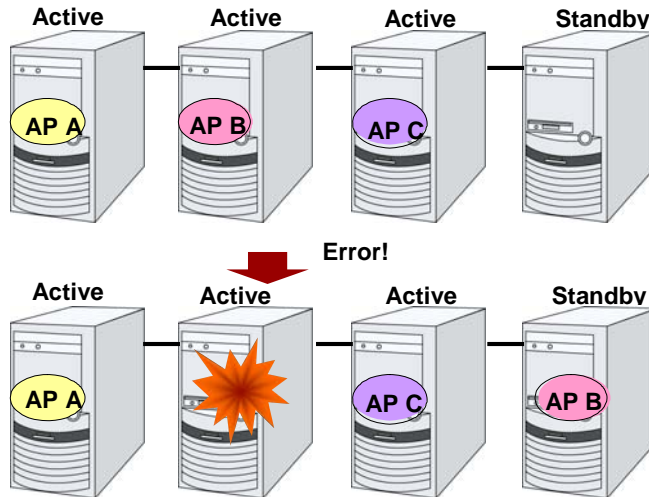


Figure 2-7 Node to Node configuration

Hardware configuration of the shared disk type cluster

The hardware configuration of the shared disk in ExpressCluster is described below. In general, the following is used for communication between the servers in a cluster system:

- ◆ Two NIC cards (one for external communication, one for ExpressCluster)
- ◆ COM port connected by RS232C cross cable
- ◆ Specific space of a shared disk

SCSI or FibreChannel can be used for communication interface to a shared disk; however, recently FibreChannel is more commonly used.

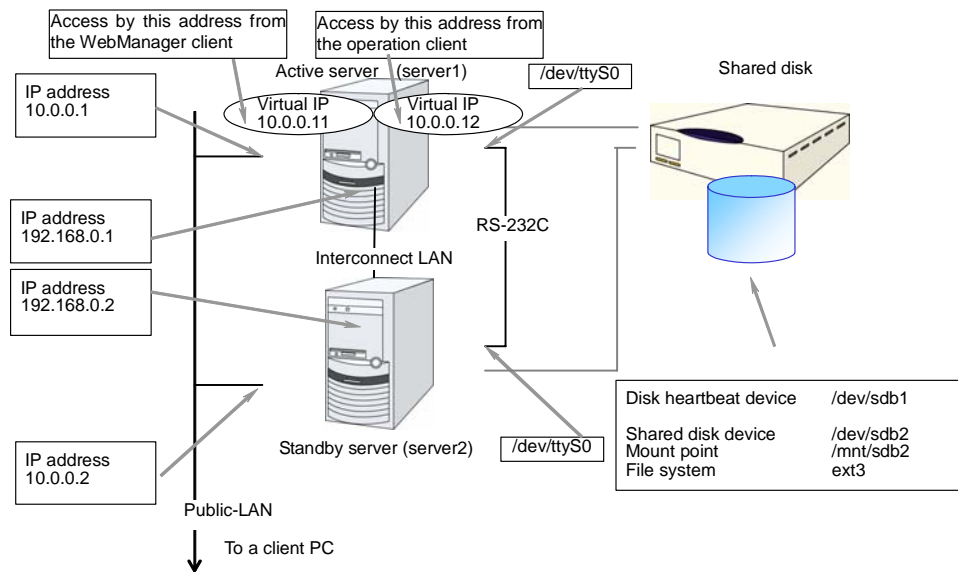


Figure 2-8 Sample of cluster environment when a shared disk is used

Hardware configuration of the mirror disk type cluster

The hardware configuration of the mirror disk in ExpressCluster is described below.

Unlike the shared disk type, a network to copy the mirror disk data is necessary. In general, a network is used with NIC for internal communication in ExpressCluster.

Mirror disks need to be separated from the operating system; however, they do not depend on a connection interface (IDE or SCSI.)

Figure 2-9 Sample of cluster environment when mirror disks are used (when allocating cluster partition and data partition to the disk where OS is installed):

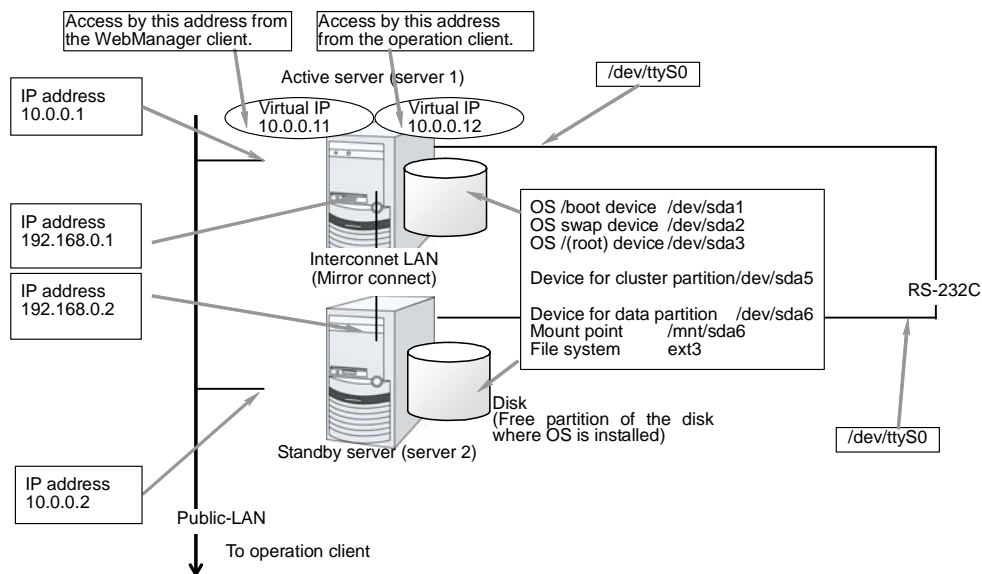
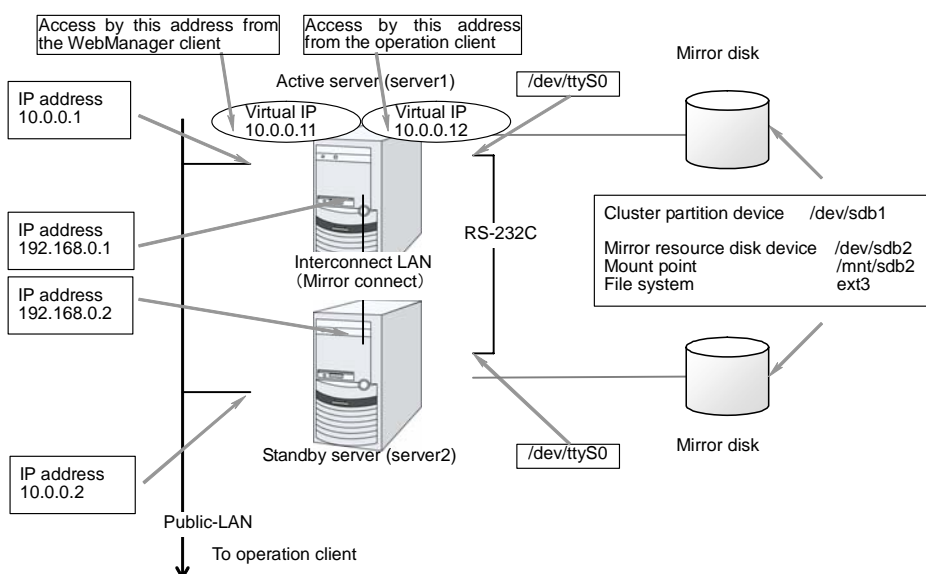


Figure 2-10 Sample of cluster environment when mirror disks are used (when disks for cluster partition and data partition are prepared):



Hardware configuration of the hybrid disk type cluster

The hardware configuration of the hybrid disk in ExpressCluster is described below.

Unlike the shared disk type, a network to copy the data is necessary. In general, NIC for internal communication in ExpressCluster is used to meet this purpose.

Disks do not depend on a connection interface (IDE or SCSI).

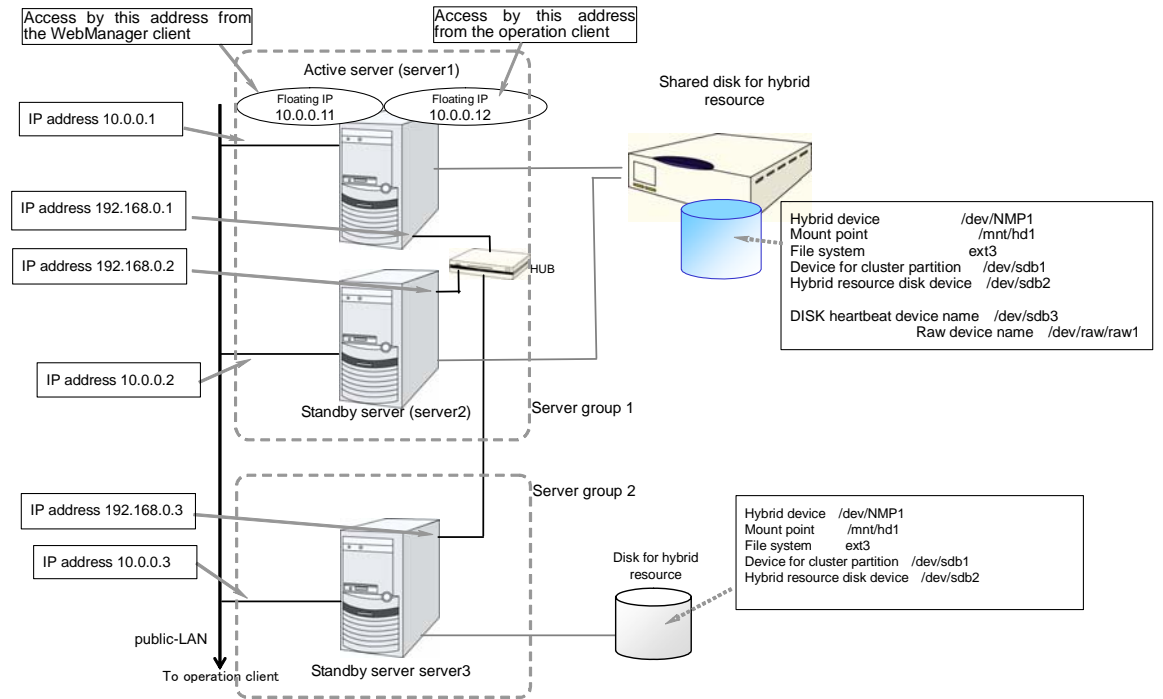


Figure 2-11: Sample of cluster environment where hybrid disks are used (two servers use a shared disk and the third server's general disk are used for mirroring)

What is cluster object?

In ExpressCluster, the various resources are managed as the following groups:

Cluster object

Configuration unit of a cluster.

Server object

Indicates the physical server and belongs to the cluster object.

Server group object

Groups the servers and belongs to the cluster object.

Heartbeat resource object

Indicates the network part of the physical server and belongs to the server object.

Network partition resolution resource object

Indicates the network partition resolution mechanism and belongs to the server object.

Group object

Indicates a virtual server and belongs to the cluster object.

Group resource object

Indicates resources (network, disk) of the virtual server and belongs to the group object.

Monitor resource object

Indicates monitoring mechanism and belongs to the cluster object.

What is a resource?

In ExpressCluster, a group used for monitoring the target is called “resources.” There are four types of resources and are managed separately. Having resources allows distinguishing what is monitoring and what is being monitored more clearly. It also makes building a cluster and handling an error easy. The resources can be divided into heartbeat resources, network partition resolution resources, group resources, and monitor resources.

Heartbeat resources

Heartbeat resources are used for verifying whether the other server is working properly between servers. The following heartbeat resources are currently supported:

LAN heartbeat resource

Uses Ethernet for communication.

Kernel mode LAN heartbeat resource

Uses Ethernet for communication.

COM heartbeat resource

Uses RS232C (COM) for communication.

Disk heartbeat resource

Uses a specific partition (cluster partition for disk heartbeat) on the shared disk for communication. It can be used only on a shared disk configuration.

Network partition resolution resources

The resource used for solving the network partition is shown below:

PING network partition resolution resource

This is a network partition resolution resource by the PING method.

Group resources

A group resource constitutes a unit when a failover occurs. The following group resources are currently supported:

Floating IP resource (fip)

Provides a virtual IP address. A client can access virtual IP address the same way as the regular IP address.

EXEC resource (exec)

Provides a mechanism for starting and stopping the applications such as DB and httpd.

Disk resource (disk)

Provides a specified partition on the shared disk. It can be used only on a shared disk configuration.

Mirror disk resource (md)

Provides a specified partition on the mirror disk. It can be used only on a mirror disk configuration.

Hybrid disk resource (hd)

Provides a specified partition on a shared disk or a disk. It can be used only for hybrid configuration.

Volume manager resource (volmgr)

Handles multiple storage devices and disks as a single logical disk.

NAS resource (nas)

Connect to the shared resources on NAS server. Note that it is not a resource that the cluster server behaves as NAS server.

Virtual IP resource (vip)

Provides a virtual IP address. This can be accessed from a client in the same way as a general IP address. This can be used in the remote cluster configuration among different network addresses.

VM resource (vm)

Starts, stops, or migrates the virtual machine.

Dynamic DNS resource (ddns)

Registers the virtual host name and the IP address of the active server to the dynamic DNS server.

Monitor resources

A monitor resource monitors a cluster system. The following monitor resources are currently supported:

IP monitor resource (ipw)

Provides a monitoring mechanism of an external IP address.

Disk monitor resource (diskw)

Provides a monitoring mechanism of the disk. It also monitors the shared disk.

Mirror disk monitor resource (mdw)

Provides a monitoring mechanism of the mirroring disks.

Mirror disk connect monitor resource (mdnw)

Provides a monitoring mechanism of the mirror disk connect.

Hybrid disk monitor resource (hdw)

Provides a monitoring mechanism of the hybrid disk.

Hybrid disk connect monitor resource (hdnw)

Provides a monitoring mechanism of the hybrid disk connect.

PID monitor resource (pidw)

Provides a monitoring mechanism to check whether a process started up by exec resource is active or not.

User mode monitor resource (userw)

Provides a monitoring mechanism for a stalling problem in the user space.

NIC Link Up/Down monitor resource (miw)

Provides a monitoring mechanism for link status of LAN cable.

Volume manager monitor resource (volmgrw)

Provides a monitoring mechanism for multiple storage devices and disks.

Multi target monitor resource (mtw)

Provides a status with multiple monitor resources.

Virtual IP monitor resource (vipw)

Provides a mechanism for sending RIP packets of a virtual IP resource.

ARP monitor resource (arpw)

Provides a mechanism for sending ARP packets of a floating IP resource or a virtual IP resource.

Custom monitor resource (genw)

Provides a monitoring mechanism to monitor the system by the operation result of commands or scripts which perform monitoring, if any.

DB2 monitor resource (db2w)

Provides a monitoring mechanism for IBM DB2 database.

ftp monitor resource (ftpw)

Provides a monitoring mechanism for FTP server.

http monitor resource (httpw)

Provides a monitoring mechanism for HTTP server.

imap4 monitor resource (imap4w)

Provides a monitoring mechanism for IMAP4 server.

MySQL monitor resource (mysqlw)

Provides a monitoring mechanism for MySQL database.

nfs monitor resource (nfsw)

Provides a monitoring mechanism for nfs file server.

Oracle monitor resource (oraclew)

Provides a monitoring mechanism for Oracle database.

OracleAS monitor resource (oracleasw)

Provides a monitoring mechanism for Oracle application.

pop3 monitor resource (pop3w)

Provides a monitoring mechanism for POP3 server.

PostgreSQL monitor resource (psqlw)

Provides a monitoring mechanism for PostgreSQL database.

samba monitor resource (sambaw)

Provides a monitoring mechanism for samba file server.

smtp monitor resource (smtpw)

Provides a monitoring mechanism for SMTP server.

Sybase monitor resource (sybasew)

Provides a monitoring mechanism for Sybase database.

Tuxedo monitor resource (tuxw)

Provides a monitoring mechanism for Tuxedo application server.

Websphere monitor resource (wasw)

Provides a monitoring mechanism for Websphere application server.

Weblogic monitor resource (wls)

Provides a monitoring mechanism for Weblogic application server.

WebOTX monitor resource (otxsw)

Provides a monitoring mechanism for WebOTX application server.

VM monitor resource (vmw)

Checks whether the virtual machine is alive.

Message receive monitor resource (mrw)

Specifies the action to take when an error message is received and how the message is displayed on the WebManager.

Dynamic DNS monitor resource (ddnsw)

Periodically registers the virtual host name and the IP address of the active server to the dynamic DNS server.

Getting started with ExpressCluster

Refer to the following guides when building a cluster system with ExpressCluster:

Latest information

Refer to Section II, “Installing ExpressCluster” in this guide.

Designing a cluster system

Refer to Section I, “Configuring a cluster system” in the *Installation and Configuration Guide* and Section II, “Resource details” in the *Reference Guide*.

Configuring a cluster system

Refer to the *Installation and Configuration Guide*. When using an optional monitoring command, refer to the *Administrator’s Guide* that is available for each target monitoring application.

Troubleshooting the problem

Refer to Section III, “Maintenance information” in the *Reference Guide*.

Section II Installing ExpressCluster

This section provides the latest information on the ExpressCluster. The latest information on the supported hardware and software is described in detail. Topics such as restrictions, known problems, and how to troubleshoot the problem are covered.

- Chapter 3 Installation requirements for ExpressCluster
- Chapter 4 Latest version information
- Chapter 5 Notes and Restrictions
- Chapter 6 Upgrading ExpressCluster

Chapter 3 Installation requirements for ExpressCluster

This chapter provides information on system requirements for ExpressCluster.
This chapter covers:

- Hardware..... 52
- Software..... 54
- System requirements for the Builder 62
- System requirements for the WebManager 63
- System requirements for the Integrated WebManager 64

Hardware

ExpressCluster operates on the following server architectures:

- ◆ IA-32
- ◆ x86_64
- ◆ ppc64 (Replicator, Replicator DR, Agents except DB Agent are not supported)

General server requirements

Required specifications for ExpressCluster Server are the following:

- ◆ RS-232C port 1 port (not necessary when configuring a cluster with 3 or more nodes)
- ◆ Ethernet port 2 or more ports
- ◆ Shared disk
- ◆ Mirror disk or empty partition for mirror
- ◆ CD-ROM drive

When using the off-line Builder upon constructing and changing the existing configuration, one of the following is required for communication between the off-line Builder and servers:

- ◆ Removable media (for example, floppy disk drive or USB flash drive)
- ◆ A machine to operate the off-line Builder and a way to share files

Supported disk interfaces

Disk types that are supported as mirror disks or hybrid disk (non-shared disk) of Replicator DR are as follows:

Disk type	Host side driver	Remarks
IDE	ide	Supported up to 120GB
SCSI	aic7xxx	
SCSI	aic79xx	
SCSI	sym53c8xx	
SCSI	mptbase,mptscsih	
SCSI	mptsas	
RAID	Megaraid (SCSI type)	
RAID	megaraid (IDE type)	Supported up to 275GB
S-ATA	sata-nv	Supported up to 80GB
S-ATA	ata-piix	Supported up to 120GB

Supported network interfaces

The following are the network boards that are supported as a mirror disk connect for the mirror disk and hybrid disk of the Replicator and the Replicator DR:

Chip	Driver
Intel 82540EM	e1000
Intel 82544EI	
Intel 82546EB	
Intel 82546GB	
Intel 82573L	
Intel 80003ES2LAN	
Intel 631xESB/632xESB	
Broadcom BCM5701	bcm5700
Broadcom BCM5703	
Broadcom BCM5721	
Broadcom BCM5721	tg3

Only typical examples are listed above and other products can also be used.

Software

System requirements for ExpressCluster Server

Supported distributions and kernel versions

The environment where ExpressCluster Server can operate depends on kernel module versions because there are kernel modules unique to ExpressCluster. Kernel versions that provide the complying kernel module are listed below.

ExpressCluster Server only runs on the kernel versions listed below.

IA-32

Distribution	Kernel version	Replicator Replicator DR support	Run clpka and clpkhb support	Express Cluster Version	Rem arks
Turbolinux 11 Server (SP1)	2.6.23-10 smp64G-2.6.23-10	Yes	Yes	3.0.0-1~	
	2.6.23-12 smp64G-2.6.23-12	Yes	Yes	3.0.0-1~	
Turbolinux Appliance Server 3.0 (SP1)	2.6.23-10 smp64G-2.6.23-10	Yes	Yes	3.0.0-1~	
	2.6.23-12 smp64G-2.6.23-12	Yes	Yes	3.0.0-1~	
Red Hat Enterprise Linux 5 (update4)	2.6.18-164.el5 PAE-2.6.18-164.el5 xen-2.6.18-164.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.6.1.el5 PAE-2.6.18-164.6.1.el5 xen-2.6.18-164.6.1.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.9.1.el5 PAE-2.6.18-164.9.1.el5 xen-2.6.18-164.9.1.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.11.1.el5 PAE-2.6.18-164.11.1.el5 xen-2.6.18-164.11.1.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.15.1.el5 PAE-2.6.18-164.15.1.el5 xen-2.6.18-164.15.1.el5	Yes	Yes	3.0.0-1~	
Red Hat Enterprise Linux 5 (update5)	2.6.18-194.el5 PAE-2.6.18-194.el5 xen-2.6.18-194.el5	Yes	Yes	3.0.0-1~	
	2.6.18-194.8.1.el5 PAE-2.6.18-194.8.1.el5 xen-2.6.18-194.8.1.el5	Yes	Yes	3.0.0-1~	
Asianux Server 3 (SP2)	2.6.18-128.7AXS3 2.6.18-128.7AXS3PAE 2.6.18-128.7AXS3xen	Yes	Yes	3.0.0-1~	
Asianux Server 3 (SP3)	2.6.18-194.1.AXS3 2.6.18-194.1.AXS3PAE 2.6.18-194.1.AXS3xen	Yes	Yes	3.0.0-1~	

Distribution	Kernel version	Replicator Replicator DR support	Run clpka and clpkhb support	Express Cluster Version	Rem arks
	2.6.18-194.2.AXS3 2.6.18-194.2.AXS3PAE 2.6.18-194.2.AXS3xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 10 (SP2)	2.6.16.60-0.21-default 2.6.16.60-0.21-smp 2.6.16.60-0.21-bigsmpt 2.6.16.60-0.21-xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 10 (SP3)	2.6.16.60-0.54.5-default 2.6.16.60-0.54.5-smp 2.6.16.60-0.54.5-bigsmpt 2.6.16.60-0.54.5-xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 11	2.6.27.19-5.1-default 2.6.27.19-5.1-pae 2.6.27.19-5.1-xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 11 (SP1)	2.6.27.12-0.7-default 2.6.27.12-0.7-pae 2.6.27.12-0.7-xen	No	Yes	3.0.0-1~	
XenServer 5.5	2.6.18-128.1.6.el5.xs5.5.0.5 05.1024xen	No	Yes	3.0.0-1~	

x86_64

Distribution	Kernel version	Replicator DR support	Run clpka and clpkhb support	Express Cluster Version	Remarks
Turbolinux 11 Server (SP1)	2.6.23-10	Yes	Yes	3.0.0-1~	
	2.6.23-12	Yes	Yes	3.0.0-1~	
Turbolinux Appliance Server 3.0 (SP1)	2.6.23-10	Yes	Yes	3.0.0-1~	
	2.6.23-12	Yes	Yes	3.0.0-1~	
Red Hat Enterprise Linux 5 (update4)	2.6.18-164.el5 xen-2.6.18-164.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.6.1.el5 xen-2.6.18-164.6.1.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.9.1.el5 xen-2.6.18-164.9.1.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.11.1.el5 xen-2.6.18-164.11.1.el5	Yes	Yes	3.0.0-1~	
	2.6.18-164.15.1.el5 xen-2.6.18-164.15.1.el5	Yes	Yes	3.0.0-1~	
Red Hat Enterprise Linux 5 (update5)	2.6.18-194.el5 xen-2.6.18-194.el5	Yes	Yes	3.0.0-1~	
	2.6.18-194.8.1.el5 xen-2.6.18-194.8.1.el5	Yes	Yes	3.0.0-1~	
Asianux Server 3 (SP2)	2.6.18-128.7AXS3 2.6.18-128.7AXS3xen	Yes	Yes	3.0.0-1~	
Asianux Server 3 (SP3)	2.6.18-194.1.AXS3 2.6.18-194.1.AXS3xen	Yes	Yes	3.0.0-1~	
	2.6.18-194.2.AXS3 2.6.18-194.2.AXS3xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 10 (SP2)	2.6.16.60-0.21-default 2.6.16.60-0.21-smp 2.6.16.60-0.21-xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 10 (SP3)	2.6.16.60-0.54.5-default 2.6.16.60-0.54.5-smp 2.6.16.60-0.54.5-xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 11	2.6.27.19-5.1-default 2.6.27.19-5.1-xen	Yes	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 11 (SP1)	2.6.27.19-5.1-default 2.6.27.19-5.1-xen	No	Yes	3.0.0-1~	
Oracle Enterprise Linux 5 (5.5)	2.6.18-194.el5	Yes	Yes	3.0.0-1~	
VMware ESX 4.0	2.6.18-128.ESX	No	Yes	3.0.0-1~	
VMware ESX 4.1	2.6.18-164.ESX	No	Yes	3.0.0-1~	

ppc64

Distribution	Kernel version	Replicator Replicator DR support	Run clpka and clpkhb support	Express Cluster Version	Rem arks
Red Hat Enterprise Linux 5 (update4)	2.6.18-164.el5	No	Yes	3.0.0-1~	
Red Hat Enterprise Linux 5 (update5)	2.6.18-194.el5	No	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 10 (SP2)	2.6.16.60-0.21-default	No	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 10 (SP3)	2.6.16.60-0.54.5-default	No	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 11	2.6.27.19-5.1-default	No	Yes	3.0.0-1~	
Novell SUSE LINUX Enterprise Server 11 (SP1)	2.6.32.12-0.7-default	No	Yes	3.0.0-1~	

Applications supported by monitoring options

Version information of the applications to be monitored by monitor resources is described below.

For the support information on the monitoring options of command type (that are registered as script resources at setup), which is provided on ExpressCluster 3.0.x-x, see the administrator's guide of each option.

IA32

Monitor resource	Monitored application	ExpressCluster version	Remarks
Oracle monitor	Oracle Database 10g Release 2 (10.2)	3.0.0-1~	
	Oracle Database 11g Release 1 (11.1)	3.0.0-1~	
	Oracle Database 11g Release 2 (11.2)	3.0.0-1~	
DB2 monitor	DB2 V9.5	3.0.0-1~	
	DB2 V9.7	3.0.0-1~	
PostgreSQL monitor	PostgreSQL 8.1	3.0.0-1~	
	PostgreSQL 8.2	3.0.0-1~	
	PostgreSQL 8.3	3.0.0-1~	
	PostgreSQL 8.4	3.0.0-1~	
	PowerGres on Linux 6.0	3.0.0-1~	
	PowerGres on Linux 7.0	3.0.0-1~	
MySQL monitor	MySQL 5.0	3.0.0-1~	
	MySQL 5.1	3.0.0-1~	
	MySQL 5.4	3.0.0-1~	
Sybase monitor	Sybase 15.0	3.0.0-1~	
Samba monitor	Samba 3.0	3.0.0-1~	
	Samba 3.2	3.0.0-1~	
	Samba 3.3	3.0.0-1~	
	Samba 3.4	3.0.0-1~	
nfs monitor	No specified version	3.0.0-1~	
HTTP monitor	No specified version	3.0.0-1~	
SMTP monitor	No specified version	3.0.0-1~	
pop3 monitor	No specified version	3.0.0-1~	
imap4 monitor	No specified version	3.0.0-1~	
ftp monitor	No specified version	3.0.0-1~	
Tuxedo monitor	Tuxedo 10g R3	3.0.0-1~	
	Tuxedo 11g R1	3.0.0-1~	
OracleAS monitor	Oracle Application Server 10g Release 3 (10.1.3)	3.0.0-1~	

	Oracle Application Server 10g Release 4 (10.1.4)	3.0.0-1~	
Weblogic monitor	WebLogic Server 10g R3	3.0.0-1~	
	WebLogic Server 11g R1	3.0.0-1~	
Websphere monitor	WebSphere 6.1	3.0.0-1~	
	WebSphere 7.0	3.0.0-1~	
WebOTX monitor	WebOTX V7.1	3.0.0-1~	
	WebOTX V8.0	3.0.0-1~	
	WebOTX V8.1	3.0.0-1~	
	WebOTX V8.2	3.0.0-1~	

x86_64

Monitor resource	Monitored application	ExpressCluster version	Remarks
Oracle monitor	Oracle Database 10g Release 2 (10.2)	3.0.0-1~	
	Oracle Database 11g Release 1 (11.1)	3.0.0-1~	
	Oracle Database 11g Release 2 (11.2)	3.0.0-1~	
DB2 monitor	DB2 V9.5	3.0.0-1~	
	DB2 V9.7	3.0.0-1~	
PostgreSQL monitor	PostgreSQL 8.1	3.0.0-1~	
	PostgreSQL 8.2	3.0.0-1~	
	PostgreSQL 8.3	3.0.0-1~	
	PostgreSQL 8.4	3.0.0-1~	
	PowerGres on Linux 6.0	3.0.0-1~	
	PowerGres on Linux 7.0	3.0.0-1~	
	PowerGres on Linux 7.1	3.0.0-1~	
	PowerGres Plus V5.0	3.0.0-1~	
MySQL monitor	MySQL 5.0	3.0.0-1~	
	MySQL 5.1	3.0.0-1~	
	MySQL 5.4	3.0.0-1~	
Sybase monitor	Sybase 15.0	3.0.0-1~	
Samba monitor	Samba 3.0	3.0.0-1~	
	Samba 3.2	3.0.0-1~	
	Samba 3.3	3.0.0-1~	
	Samba 3.4	3.0.0-1~	
NFS monitor	No specified version	3.0.0-1~	
HTTP monitor	No specified version	3.0.0-1~	
SMTP monitor	No specified version	3.0.0-1~	
pop3 monitor	No specified version	3.0.0-1~	

imap4 monitor	No specified version	3.0.0-1~	
ftp monitor	No specified version	3.0.0-1~	
Tuxedo monitor	Tuxedo 10g R3	3.0.0-1~	
	Tuxedo 11g R1	3.0.0-1~	
OracleAS monitor	Oracle Application Server 10g Release 3 (10.1.3)	3.0.0-1~	
	Oracle Application Server 10g Release 4 (10.1.4)	3.0.0-1~	
Weblogic monitor	WebLogic Server 10g R3	3.0.0-1~	
	WebLogic Server 11g R3	3.0.0-1~	
Websphere monitor	WebSphere 6.1	3.0.0-1~	
	WebSphere 7.0	3.0.0-1~	
WebOTX monitor	WebOTX V7.1	3.0.0-1~	
	WebOTX V8.0	3.0.0-1~	
	WebOTX V8.1	3.0.0-1~	
	WebOTX V8.2	3.0.0-1~	

ppc64

Monitor resource	Monitored application	ExpressCluster version	Remarks
Oracle monitor	Oracle Database 10g Release 2 (10.2)	3.0.0-1~	
DB2 monitor	DB2 V9.7	3.0.0-1~	
PostgreSQL monitor	PostgreSQL 8.4	3.0.0-1~	

Operation Environment of VM resources

The followings are the version information of the virtual machines on which VM resources operation are verified.

Virtual Machine	Version	ExpressCluster version	Remarks
vSphere	4.0 update1	3.0.0-1~	x86_64
	4.0 update2	3.0.0-1~	x86_64
	4.1	3.0.0-1~	x86_64
XenServer	5.5	3.0.0-1~	IA32
KVM	Redhat Enterprise Linux 5.5	3.0.0-1~	x86_64

Required memory and disk size

	Required memory size		Required disk size	
	User mode	Kernel mode	Right after installation	Max. during operation
IA-32	64MB	When the synchronization mode is used:	140MB	1.2GB
x86_64	64MB	(number of request queues x I/O size) + (2MB x number of mirror disk resources and hybrid disk resources) When the asynchronous mode is used: (number of request queues x I/O size) + ((2MB + (number of asynchronous queues)) x number of mirror disk resources and hybrid disk resources)	140MB	1.2GB
ppc64	64MB	-	24MB	1.1GB

Note:

The I/O size is 128 KB for the vxfs file system and 4KB for file systems other than it.

For the setting value of the number of request queues and asynchronization queues, see “Understanding mirror disk resources” in the *Reference Guide*.

System requirements for the Builder

Supported operating systems and browsers

Refer to the website, <http://www.nec.com/expresscluster/>, for the latest information. Currently supported operating systems and browsers are the following:

Operating system	Browser	Language
Microsoft Windows® XP SP3 or later (IA32)	IE7	English/Japanese/Chinese
	IE8	English/Japanese/Chinese
Microsoft Windows Vista™ (IA32)	IE7	English/Japanese/Chinese
	IE8	English/Japanese/Chinese
Microsoft Windows® 7 (IA32)	IE8	English/Japanese/Chinese
Microsoft Windows Server 2008 (IA32)	IE7	English/Japanese/Chinese
Novell SUSE LINUX Enterprise Server 10 (IA32)	Firefox 2.0.0.2	English/Japanese/Chinese
Red Hat Enterprise Linux 5 update5 (IA32)	Firefox 3.0.18	English/Japanese/Chinese
Asianux Server 3 (IA32)	Firefox 1.5.0.12	English/Japanese/Chinese
	Konqueror3.5.5	English/Japanese/Chinese
Turbolinux 11 Server (IA32)	Firefox 2.0.0.8	English/Japanese/Chinese

Note:

The Builder does not operate on x86_64, ppc64 machines. Use 32-bit machine when constructing and changing a cluster configuration.

Java runtime environment

Required:

Sun Microsystems, Java™ Runtime Environment, Version 6.0 Update 21 (1.6.0_21) or later

Required memory and disk size

Required memory size: 32MB or more

Required disk size: 5MB (excluding the size required for Java runtime environment)

Supported ExpressCluster versions

Offline Builder version	ExpressCluster X rpm version
3.0.0-1	3.0.0-1

Note:

When you use the Offline Builder and the ExpressCluster rpm, a combination of their versions should be the one shown above. The Builder may not operate properly if they are used in a different combination.

System requirements for the WebManager

Supported operating systems and browsers

Refer to the website, <http://www.nec.com/expresscluster/>, for the latest information. Currently the following operating systems and browsers are supported:

Operating system	Browser	Language
Microsoft Windows® SP3(IA32)	IE6 SP2	English/Japanese/Chinese
	IE7	English/Japanese/Chinese
	IE8	English/Japanese/Chinese
Microsoft Windows Vista® (IA32)	IE7	English/Japanese/Chinese
	IE8	English/Japanese/Chinese
Microsoft Windows® 7 (IA32)	IE8	English/Japanese/Chinese
Microsoft Windows Server 2008 (IA32)	IE7	English/Japanese/Chinese
Novell SUSE LINUX Enterprise Server 10 (IA32)	Firefox 2.0.0.2	English/Japanese/Chinese
Red Hat Enterprise Linux 5 update1 (IA32)	Firefox 1.5.0.12	English/Japanese/Chinese
Asianux Server 3 (IA32)	Firefox 1.5.0.12	English/Japanese/Chinese
	Konqueror3.5.5	English/Japanese/Chinese
Turbolinux 11 Server (IA32)	Firefox 2.0.0.8	English/Japanese/Chinese

Note:

The ExpressCluster X WebManager does not run on x86_64, and PPC 64 machines. Use a 32-bit OS when operating a cluster on Linux machine.

Java runtime environment

Required:

Sun Microsystems, Java™ Runtime Environment, Version 6.0 Update 21 (1.6.0_21) or later

Required memory and disk size

Required memory size: 40MB or more

Required disk size: 600KB (excluding the size required for Java runtime environment)

System requirements for the Integrated WebManager

This section explains system requirements to operate the Integrated WebManager. Refer to “Integrated WebManager Administrator’s Guide” for the Java application version Integrated WebManager.

Supported operating systems and browsers

Currently the following operating systems and browsers are supported:

Operating system	Browser	Language
Windows® XP SP2	IE6 SP2	English/Japanese/Chinese
Microsoft Windows Vista™	IE7	English/Japanese/Chinese
Microsoft Windows® 7 (IA32)	IE8	English/Japanese/Chinese
Microsoft Windows Server 2003 SP1	IE6 SP1	English/Japanese/Chinese
Microsoft Windows Server 2003 R2	IE6 SP1	English/Japanese/Chinese
Microsoft Windows Server 2008	IE7	English/Japanese/Chinese

Note:

To use the WebManager on x86_64 machines, it is necessary to use the 32-bit Java Runtime. For the latest information, refer to ExpressCluster website, <http://www.nec.co.jp/pfsoft/clusterpro/clp/overseas.html>

Java runtime environment

Required:

Sun Microsystems, Java(TM) Runtime Environment, Version 6.0 Update 21 (1.6.0_21) or later (noncompliant to Version 6.0)

Required memory size and disk size

Required memory size: 40MB or more

Required disk size: 300KB or more (excluding the size required for Java runtime environment)

Chapter 4 Latest version information

This chapter provides the latest information on ExpressCluster.
This chapter covers:

- Correspondence list of ExpressCluster and a manual 66
- Enhanced functions 67
- Corrected information 69

Correspondence list of ExpressCluster and a manual

This book has explained on the assumption that ExpressCluster of the following version. Be careful of the number of versions of the version of ExpressCluster, and a manual.

ExpressCluster Version	Manual	Manual Version	Remarks
3.0.0-1	Installation and Configuration Guide	First Edition	
	Getting Started Guide	First Edition	
	Reference Guide	First Edition	
	Integrated WebManager Administrator's Guide	First Edition	

Enhanced functions

Upgrade has been performed on the following minor versions.

Number	Version (in detail)	Upgraded section
1	3.0.0-1	The WebManager and Builder can now be used from the same browser window.
2	3.0.0-1	The cluster configuration wizard has been upgraded.
3	3.0.0-1	Some settings can now be automatically acquired in the cluster configuration wizard.
4	3.0.0-1	The Integrated WebManager can now be used from a browser.
5	3.0.0-1	A function has been implemented to check settings when uploading configuration data.
6	3.0.0-1	ExpressCluster can now automatically select the failover destination when an error occurs.
7	3.0.0-1	A function has been implemented to control failovers across server groups.
8	3.0.0-1	All Groups can now be selected as the failover target when an error is detected.
9	3.0.0-1	The start wait time can now be skipped.
10	3.0.0-1	ExpressCluster can now manage external errors.
11	3.0.0-1	Dump information can now be acquired when the target monitoring application times out.
12	3.0.0-1	Detailed information about an Oracle database can now be acquired if an error is detected while monitoring it.
13	3.0.0-1	Mirror data can now be compressed for transfer during asynchronous mirroring.
14	3.0.0-1	Whole mirror synchronization has been accelerated.
15	3.0.0-1	A function has been implemented to register a virtual host name to the dynamic DNS server.
16	3.0.0-1	A guest OS can now be handled as a resource when the host OS of vSphere, XenServer, or kvm is clustered.
17	3.0.0-1	A function has been implemented to automatically follow a guest OS in the virtualization infrastructure if it is moved by software other than ExpressCluster.
18	3.0.0-1	vMotion can now be executed at error detection or during operation if the vSphere host OS is clustered.
19	3.0.0-1	The Logical Volume Manager (LVM) can now be controlled.
20	3.0.0-1	Disk settings have been consolidated.

21	3.0.0-1	Additional OSs are now supported.
22	3.0.0-1	Additional applications are now supported.
23	3.0.0-1	Additional network warning lights are now supported.

Corrected information

There is no corrected information because this is the first edition.

Chapter 5 Notes and Restrictions

This chapter provides information on known problems and how to troubleshoot the problems.

This chapter covers:

- Designing a system configuration..... 72
- Before installing ExpressCluster 83
- Notes when creating ExpressCluster configuration data 88
- After start operating ExpressCluster..... 94

Designing a system configuration

Hardware selection, option products license arrangement, system configuration, and shared disk configuration are introduced in this section.

Function list and necessary license

The following option products are necessary as many as the number of servers.

Necessary function	Necessary license
Mirror disk resource	CLUSTERPRO X Replicator 3.0 *1
Hybrid disk resource	CLUSTERPRO X Replicator DR 3.0 *2
Oracle monitor resource	CLUSTERPRO X Database Agent 3.0
DB2 monitor resource	CLUSTERPRO X Database Agent 3.0
PostgreSQL monitor resource	CLUSTERPRO X Database Agent 3.0
MySQL monitor resource	CLUSTERPRO X Database Agent 3.0
samba monitor resource	CLUSTERPRO X File Server Agent 3.0
nfs monitor resource	CLUSTERPRO X File Server Agent 3.0
http monitor resource	CLUSTERPRO X Internet Server Agent 3.0
smtp monitor resource	CLUSTERPRO X Internet Server Agent 3.0
pop3 monitor resource	CLUSTERPRO X Internet Server Agent 3.0
imap4 monitor resource	CLUSTERPRO X Internet Server Agent 3.0
ftp monitor resource	CLUSTERPRO X Internet Server Agent 3.0
Tuxedo monitor resource	CLUSTERPRO X Application Server Agent 3.0
OracleAS monitor resource	CLUSTERPRO X Application Server Agent 3.0
Weblogic monitor resource	CLUSTERPRO X Application Server Agent 3.0
Websphere monitor resource	CLUSTERPRO X Application Server Agent 3.0
WebOTX monitor resource	CLUSTERPRO X Application Server Agent 3.0
Sybase monitor resource	CLUSTERPRO X Application Server Agent 3.0
Mail report actions	CLUSTERPRO X Alert Service 3.0
Network Warning Light status	CLUSTERPRO X Alert Service 3.0

*1 When configuring data mirror form, product **Replicator** must be purchased.

*2 When configuring mirror between shared disk, product **Replicator DR** must be purchased.

Supported operating systems for the Builder and WebManager

- ◆ The Builder does not run on x86_64 and ppc64 machines. Use a 32-bit machine when configuring and changing the configuration of a cluster system.
- ◆ The WebManager does not run on x86_64 and ppc64 machines. Use a 32-bit machine when operating a cluster system.

Hardware requirements for mirror disks

- ◆ Disks to be used as a mirror resource do not support a Linux md and/or LVM stripe set, volume set, mirroring, and stripe set with parity.
- ◆ Mirror disk resource cannot be made as a target of a Linux md or LVM stripe set, volume set, mirroring, and stripe set with parity.
- ◆ Mirror partitions (data partition and cluster partition) to use a mirror resource.
- ◆ There are two ways to allocate mirror partitions:
 - Allocate a mirror partition (data partition and cluster partition) on the disk where the operating system (such as root partition and swap partition) resides.

- Reserve (or add) a disk (or LUN) not used by the operating system and allocate a mirror partition on the disk.
- ◆ Consider the following when allocating mirror partitions:
 - When maintainability and performance are important:
 - It is recommended to have a mirror disk that is not used by the OS.
 - When LUN cannot be added due to hardware RAID specification or when changing LUN configuration is difficult in hardware RAID pre-install model:
 - Allocate a mirror partition on the same disk where the operating system resides.
- ◆ When multiple mirror resources are used, it is recommended to prepare (adding) a disk per mirror resource. Allocating multiple mirror resources on the same disk may result in degraded performance and it may take a while to complete mirror recovery due to disk access performance on Linux operating system.
- ◆ Disks used for mirroring must be the same in all servers.

- Disk type

Mirror disks on both servers and disks where mirror partition is allocated should be of the same disk type

For supported disk types, see “Supported disk interfaces” on page 52.

Example

Combination	server1	server2
OK	SCSI	SCSI
OK	IDE	IDE
NG	IDE	SCSI

- ◆ Notes when the geometries of the disks used as mirror disks differ between the servers.

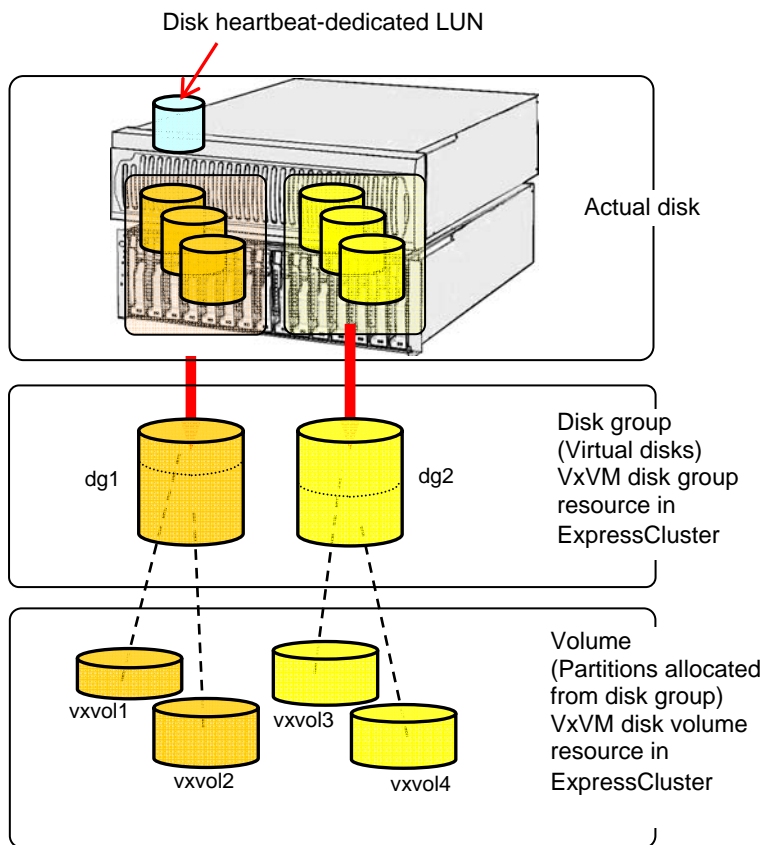
The partition size allocated by the fdisk command is aligned by the number of blocks (units) per cylinder. Allocate a data partition considering the relationship between data partition size and direction for initial mirror configuration to be as indicated below:

Source server ≤ Destination server

“Source server” refers to the server where the failover group that a mirror resource belongs has a higher priority in failover policy. “Destination server” refers to the server where the failover group that a mirror resource belongs has a lower priority in failover policy.

Hardware requirements for shared disks

- ◆ A shared disk does not support a Linux md stripe set, volume set, mirroring, and stripe set with parity.
- ◆ When a Linux LVM stripe set, volume set, mirroring, or stripe set with parity is used:
 - ExpressCluster cannot control ReadOnly/ReadWrite of the partition configured for the disk resource.
- ◆ When you use VxVM or LVM, a LUN that is not controlled by VxVM or LVM is required on a shared disk for the disk heartbeat of ExpressCluster. You should bear this in your mind when configuring LUN on the shared disk.



Hardware requirements for hybrid disks

- ◆ Disks to be used as a hybrid resource do not support a Linux md and/or LVM stripe set, volume set, mirroring, and stripe set with parity.
- ◆ Hybrid disk resource cannot be made as a target of a Linux md or LVM stripe set, volume set, mirroring, and stripe set with parity.
- ◆ Hybrid partitions (data partition and cluster partition) are required to use a hybrid resource.
- ◆ When a disk for hybrid disk is allocated in the shared disk, a partition for disk heartbeat resource between servers sharing the shared disk device is required.
- ◆ The following are the two ways to allocate partitions when a disk for hybrid disk is allocated from a disk which is not a shared disk:
 - Allocate hybrid partitions (data partition and cluster partition) on the disk where the operating system (such as root partition and swap partition) resides.
 - Reserve (or add) a disk (or LUN) not used by the operating system and allocate a hybrid partition on the disk.
- ◆ Consider the following when allocating hybrid partitions:
 - When maintainability and performance are important:
 - It is recommended to have a hybrid disk that is not used by the OS.
 - When LUN cannot be added due to hardware RAID specification or when changing LUN configuration is difficult in hardware RAID pre-install model:
 - Allocate a hybrid partition on the same disk where the operating system resides.
- ◆ When multiple hybrid resources are used, it is recommended to prepare (add) a LUN per hybrid resource. Allocating multiple hybrid resources on the same disk may result in degraded in performance and it may take a while to complete mirror recovery due to disk access performance on Linux operating system.

Type of required partition	Device for which hybrid disk resource is allocated	
	Shared disk device	Non-shared disk device
Data partition	Required	Required
Cluster partition	Required	Required
Partition for disk heart beat	Required	Not Required
Allocation on the same disk (LUN) as where the OS is	-	Possible

- ◆ Notes when the geometries of the disks used as hybrid disks differ between the servers.

Allocate a data partition considering the relationship between data partition size and direction for initial mirror configuration to be as indicated below:

Source server \leq Destination server

“Source server” refers to the server with a higher priority in failover policy in the failover group where the hybrid resource belongs. “Destination server” refers to the server with a lower priority in failover policy in the failover group where the hybrid resource belongs has.

NIC link up/down monitor resource

Some NIC boards and drivers do not support required `ioctl()`.

The propriety of a NIC Link Up/Down monitor resource of operation can be checked by the `ethtool` command which each distributor offers.

```
ethtool eth0
Settings for eth0:
    Supported ports: [ TP ]
    Supported link modes:   10baseT/Half 10baseT/Full
                           100baseT/Half 100baseT/Full
                           1000baseT/Full

    Supports auto-negotiation: Yes
    Advertised link modes:  10baseT/Half 10baseT/Full
                           100baseT/Half 100baseT/Full
                           1000baseT/Full

    Advertised auto-negotiation: Yes
    Speed: 1000Mb/s
    Duplex: Full
    Port: Twisted Pair
    PHYAD: 0
    Transceiver: internal
    Auto-negotiation: on
    Supports Wake-on: umbg
    Wake-on: g
    Current message level: 0x00000007 (7)
    Link detected: yes
```

- ◆ When the LAN cable link status ("Link detected: yes") is not displayed as the result of the `ethtool` command:
 - It is highly likely that NIC Link Up/Down monitor resource of EXPRESSCLUSTER is not operable. Use IP monitor resource instead.

- ◆ When the LAN cable link status ("Link detected: yes") is displayed as the result of the `ethtool` command:
 - In most cases NIC Link Up/Down monitor resource of ExpressCluster can be operated, but sometimes it cannot be operated.
 - Particularly in the following hardware, NIC Link Up/Down monitor resource of ExpressCluster may not be operated. Use IP monitor resource instead.
 - When hardware is installed between the actual LAN connector and NIC chip such as a blade server

To check if NIC Link Up/Down monitor resource can be used by using ExpressCluster on an actual machine, follow the steps below to check the operation.

1. Register NIC Link Up/Down monitor resource with the configuration information. Select **No Operation** for the configuration of recovery operation of NIC Link Up/Down monitor resource upon failure detection.
2. Start the cluster.
3. Check the status of NIC Link Up/Down monitor resource. If the status of NIC Link Up/Down monitor resource is abnormal while LAN cable link status is normal, NIC Link Up/Down monitor resource cannot be operated.

4. If NIC Link Up/Down monitor resource status becomes abnormal when LAN cable link status is made abnormal status (link down status), NIC Link Up/Down monitor resource cannot be operated.
If the status remains to be normal, NIC Link Up/Down monitor resource cannot be operated.

Write function of the mirror disk resource and hybrid disk resource

- ◆ A mirror disk and a hybrid disk resource write data in the disk of its own server and the disk of the remote server via network. Reading of data is done only from the disk on own server.
- ◆ Writing functions shows poor performance in mirroring when compared to writing to a single server because of the reason provided above. For a system that requires through-put as high as single server, use a shared disk.

Not outputting syslog to the mirror disk resource or the hybrid disk resource

Do not set directories or subdirectories which mounted the mirror disk resource or the hybrid disk resource as syslog output destination directories.

When the mirror disk connection is disconnected, the I/O to the mirror partition may stop until the disconnection is detected. The system may become abnormal because of the syslog output stoppage at this time.

When outputting syslog to the mirror disk resource or the hybrid disk resource is necessary, consider the followings.

- ◆ Use bonding as a way of path redundancy of the mirror disk connection.
- ◆ Adjust the user space monitoring timeout value or the mirror related timeout values.

Notes when terminating the mirror disk resource or the hybrid disk resource

- ◆ In case that processes which access to the directories, subdirectories and files which mounted the mirror disk resource or the hybrid disk resource exist, terminate the accesses to each disk resource by using ending script or other methods at deactivation of each disk resource like when shutdown or failover.
Depending on the settings of each disk resource, action at abnormality detection when unmounting (forcibly terminate processes while each disk resource is being accessed) may occur, or recovery action at deactivation failure caused by unmount failure (OS shutdown or other actions) may be executed.
- ◆ In case that a massive amount of accesses to directories, subdirectories or files which mounted the mirror disk resource or hybrid disk resource are executed, it may take much time before the cache of the file systems is written out to the disks when unmounting at disk resource deactivation.
At times like this, set the timeout interval of unmount longer enough so that the writing to the disks will successfully complete.
- ◆ For the details of this setting, see Chapter 4, "Group resource details" in Reference Guide, **Settings Tab** or **Mirror Disk Resource Tuning Properties** or **Unmount Tab** in **Details Tab** in "Understanding mirror disk resources" or "Understanding mirror disk resources".

Data consistency among multiple asynchronous mirror disks

In mirror disk or hybrid disk with asynchronous mode, writing data to the data partition of the active server is performed in the same order as the data partition of the standby server.

This writing order is guaranteed except during the initial mirror disk configuration or recovery (copy) period after suspending mirroring the disks. The data consistency among the files on the standby data partition is guaranteed.

However, the writing order is not guaranteed among multiple mirror disk resources and hybrid disk resources. For example, if a file gets older than the other and files that cannot maintain the data consistency are distributed to multiple asynchronous mirror disks, an application may not run properly when it fails over due to server failure.

For this reason, be sure to place these files on the same asynchronous mirror disk or hybrid disk.

O_DIRECT for mirror or hybrid disk resources

Do not use the O_DIRECT flag of the open() system call for mirror or hybrid disk resources. Examples include the Oracle parameter `filesystemio_options = setall`. Do not specify the disk monitor O_DIRECT mode for mirror or hybrid disk resources.

Initial mirror construction time for mirror or hybrid disk resources

The time that takes to construct the initial mirror is different from ext3 and other filesystems.

Installing operating system

Notes on parameters to be determined when installing an operating system, allocating resources, and naming rules are described in this section.

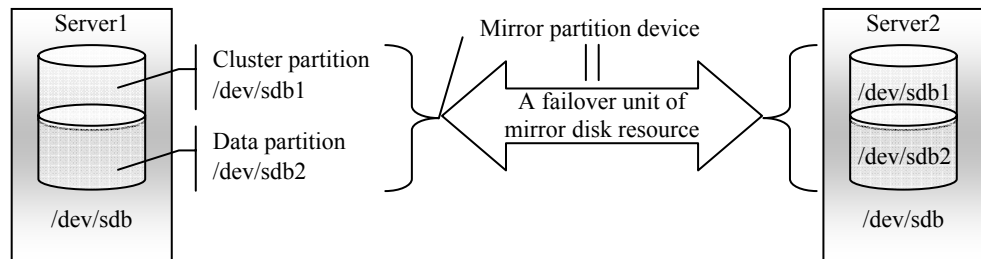
/opt/nec/clusterpro file system

It is recommended to use a file system that has journaling functions to improve tolerance for system failure.

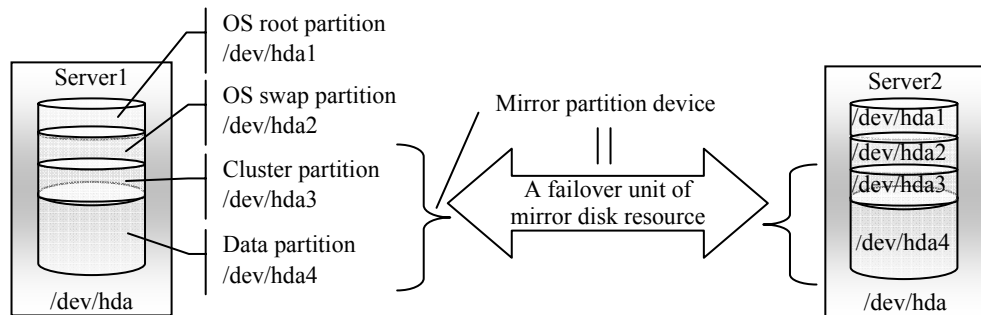
Mirror disks

◆ Disk partition

Example: When adding one SCSI disk to each of both servers and making a pair of mirrored disks:



Example: When using free space of IDE disks of both servers, where the OS is stored, and making a pair of mirrored disks:

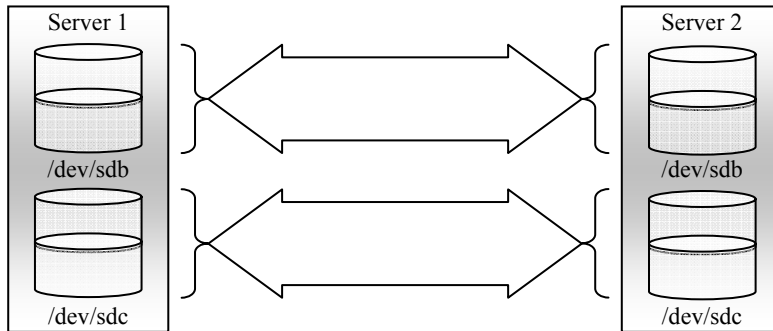


- Mirror partition device refers to cluster partition and data partition.
- Allocate cluster partition and data partition on each server as a pair.
- It is possible to allocate a mirror partition (cluster partition and data partition) on the disk where the operating system resides (such as root partition and swap partition.).
 - When maintainability and performance are important:
It is recommended to have a mirror disk that is not used by the operating system (such as root partition and swap partition.)
 - When LUN cannot be added due to hardware RAID specification: or
When changing LUN configuration is difficult in hardware RAID pre-install model:
It is possible to allocate a mirror partition (cluster partition and data partition) on the disk where the operating system resides (such as root partition and swap partition.)

◆ Disk configurations

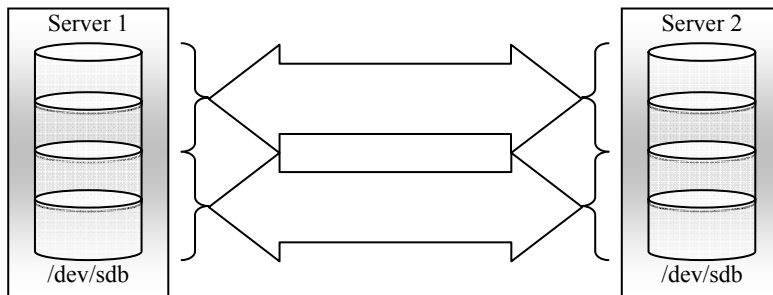
Multiple disks can be used as mirror disks on a single server. Or, you can allocate multiple mirror partitions on a single disk.

Example: When adding two SCSI disks to each of both servers and making two pairs of mirrored disks:



- Allocate two partitions, cluster partition and data partition, as a pair on each disk.
- Use of the data partition as the first disk and the cluster partition as the second disk is not permitted.

Example: When adding one SCSI disk to each of both servers and making two mirror partitions:



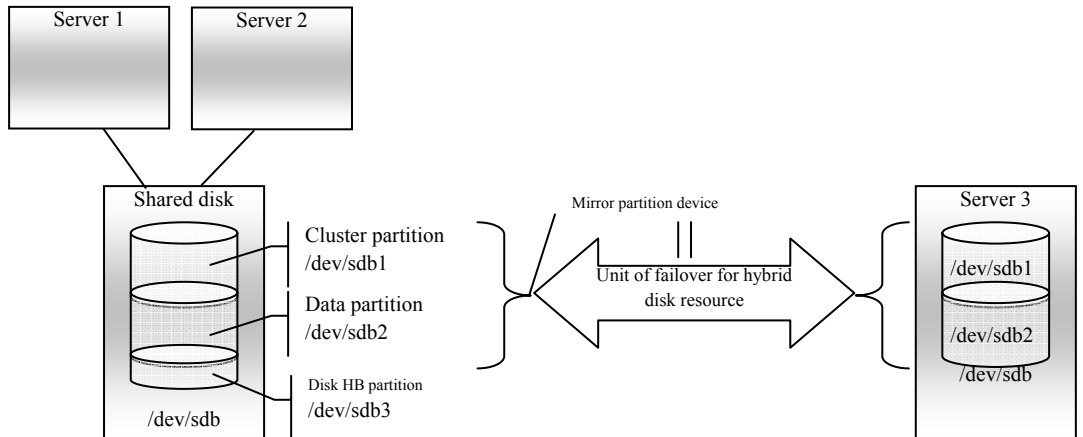
- ◆ A disk does not support a Linux md and/or LVM stripe set, volume set, mirroring, and stripe set with parity.

Hybrid disks

◆ Disk partition

Disks that are shared or not shared (server with built-in disk, external disk chassis not shared by servers etc.) can be used.

Example) When two servers use a shared disk and the third server uses a built-in disk in the server:



- Mirror partition device is a device ExpressCluster mirroring driver provides in the upper.
- Allocate cluster partition and data partition on each server as a pair.
- When a disk that is not shared (e.g. server with a built-in disk, external disk chassis that is not shared among servers) is used, it is possible to allocate mirror partitions (cluster partition and data partition) on the disk where the operating system resides (such as root partition and swap partition.).
 - When maintainability and performance are important:

It is recommended to have a mirror disk that is not used by the operating system (such as root partition and swap partition.)
 - When LUN cannot be added due to hardware RAID specification: or
When changing LUN configuration is difficult in hardware RAID pre-install model:

It is possible to allocate mirror partitions (cluster partition and data partition) on the disk where the operating system resides (such as root partition and swap partition.)
- When a hybrid disk is allocated in a shared disk device, allocate a partition for the disk heart beat resource between servers sharing the shared disk device.
- A disk does not support a Linux md and/or LVM stripe set, volume set, mirroring, and stripe set with parity.

Dependent library

◆ libxml2

Install libxml2 when installing the operating system.

Dependent driver

- ◆ softdog

This driver is necessary when softdog is used to monitor user mode monitor resource.

Configure a loadable module. Static driver cannot be used.

The major number of Mirror driver

Use mirror driver's major number 218. Do not use major number 218 for other device drivers.

The major number of Kernel mode LAN heartbeat and keepalive drivers

- ◆ Use major number 10, minor number 240 for kernel mode LAN heartbeat driver.
- ◆ Use major number 10, minor number 241 for keepalive driver.

Make sure to check that other drivers are not using major and minor numbers described above.

Partition for RAW monitoring of disk monitor resources

Allocate a partition for monitoring when setting up RAW monitoring of disk monitor resources. The partition size should be 10 MB.

SELinux settings

Configure permissive or disabled for the SELinux settings.

If you set enforcing, communication required in ExpressCluster may not be achieved.

Before installing ExpressCluster

Notes after installing an operating system, when configuring OS and disks are described in this section.

Communication port number

In ExpressCluster, the following port numbers are used. You can change the port number by using the Builder except "keepalive between mirror drivers."

Make sure not to access the following port numbers from a program other than ExpressCluster.

Configure to be able to access the port number below when setting a firewall on a server.

Server to Server			Used for		
Loopback in servers					
Server	Automatic allocation	-	Server	29001/TCP	Internal communication
Server	Automatic allocation	-	Server	29002/TCP	Data transfer
Server	Automatic allocation	-	Server	29002/UDP	Heartbeat
Server	Automatic allocation	-	Server	29003/UDP	Alert synchronization
Server	Automatic allocation	-	Server	29004/TCP	Communication between mirror agents
Server	Automatic allocation	-	Server	29006/UDP	Heartbeat (kernel mode)
Server	Automatic allocation	-	Server	XXXX/TCP	Mirror disk resource data synchronization
Server	Automatic allocation	-	Server	XXXX/TCP	Communication between mirror drivers
Server	Automatic allocation	-	Server	XXXX/TCP	Communication between mirror drivers
Server	Automatic allocation	-	Server	icmp	keepalive between mirror drivers, duplication check for FIP/VIP resource and mirror agent
Server	Automatic allocation	-	Server	XXXX/TCP	Internal log communication
WebManager to Server			Used for		
WebManager	Automatic allocation	-	Server	29003/TCP	http communication

Server connected to the Integrated WebManager to Target server					
					Used for
Server connected to the Integrated WebManager	Automatic allocation	-	Server	29003/TCP	http communication
Others					
					Used for
Server	Automatic allocation	-	Network warning light	514/TCP	Network warning light control
Server	Automatic allocation	-	Management LAN of server BMC	623/UDP	BMC control (Forced stop / Chassis lamp association)
Server	Automatic allocation	-	Monitoring target	icmp	IP monitor
Server	Automatic allocation	-	NFS server	icmp	Checking if NFS server is active by NAS resource
Server	Automatic allocation	-	Monitoring target	icmp	Monitoring target of Ping method network partition resolution resource

1. In automatic allocation, a port number not being used at a given time is allocated.
2. This is a port number used per mirror disk resource or hybrid disk resource and is set when creating mirror disk resource or hybrid disk resource. A port number 29051 is set by default. When you add a mirror disk resource or hybrid disk resource, this value is automatically incremented by 1. To change the value, click **Details** tab in the **[md] Resource Properties** or the **[hd] Resource Properties** dialog box of the Builder. For more information, refer to Chapter 4, “Group resource details” in the *Reference Guide*.
3. This is a port number used per mirror disk resource or hybrid disk resource and is set when creating mirror disk resource or hybrid disk resource. A port number 29031 is set by default. When you add a mirror disk resource or a hybrid disk resource, this value is automatically incremented by 1. To change the value, click **Details** tab in the **[md] Resource Properties** or the **[hd] Resource Properties** dialog box of the Builder. For more information, refer to Chapter 4, “Group resource details” in the *Reference Guide*.
4. This is a port number used per mirror disk resource or hybrid disk resource and is set when creating mirror disk resource or hybrid disk resource. A port number 29071 is set by default. When you add a mirror disk resource or hybrid disk resource this value is automatically incremented by 1. To change the value, click **Details** tab in the **[md] Resource Properties** or the **[hd] Resource Properties** dialog box of the Builder. For more information, refer to Chapter 4, “Group resource details” in the *Reference Guide*.
5. Select **UDP** for the **Communication Method for Internal Logs** in **Cluster Properties, Port No. (Log)** tab. Use the port number configured in Port No. Communication port is not used for the default log communication method **UNIX Domain**.

Changing the range of automatic allocation for the communication port numbers

- ◆ The range of automatic allocation for the communication port numbers managed by the OS might overlap the communication port numbers used by ExpressCluster.
- ◆ Change the OS settings to avoid duplication when the range of automatic allocation for the communication numbers managed by OS and the communication numbers used by ExpressCluster are duplicated.

Examples of checking and displaying OS setting conditions.

The range of automatic allocation for the communication port numbers depends on the distribution.

```
# cat /proc/sys/net/ipv4/ip_local_port_range
1024 65000
```

This is the condition to be assigned for the range from 1024 to 65000 when the application requests automatic allocation for the communication port numbers to the OS.

```
# cat /proc/sys/net/ipv4/ip_local_port_range
32768 61000
```

This is the condition to be assigned for the range from 32768 to 61000 when the application requests automatic allocation for the communication port numbers to the OS.

Examples of OS settings change

Add the line below to `/etc/sysctl.conf`. (When changing to the range from 30000 to 65000)

```
net.ipv4.ip_local_port_range = 30000 65000
```

Clock synchronization

In a cluster system, it is recommended to synchronize multiple server clocks regularly. Synchronize server clocks by using `ntp`.

NIC device name

Because of the `ifconfig` command specification, when the NIC device name is shortened, the length of the NIC device name which ExpressCluster can handle depends on it.

Shared disk

- ◆ When you continue using the data on the shared disk at times such as server reinstallation, do not allocate a partition or create a file system.
- ◆ The data on the shared disk gets deleted if you allocate a partition or create a file system.
- ◆ ExpressCluster controls the file systems on the shared disk. Do not include the file systems on the shared disk to `/etc/fstab` in operating system.
- ◆ See the *Installation and Configuration Guide* for steps for shared disk configuration.

Mirror disk

- ◆ Set a management partition for mirror disk resource (cluster partition) and a partition for

mirror disk resource (data partition).

- ◆ ExpressCluster controls the file systems on mirror disks. Do not set the file systems on the mirror disks to /etc/fstab in operating system.
- ◆ See the *Installation and Configuration Guide* for steps for mirror disk configuration.

Hybrid disk

- ◆ Configure the management partition (cluster partition) for hybrid disk resource and the partition used for hybrid disk resource (data partition).
- ◆ When a hybrid disk is allocated in the shared disk device, allocate the partition for the disk heart beat resource between servers sharing the shared disk device.
- ◆ ExpressCluster controls the file systems on the hybrid disk. Do not include the file systems on the hybrid disk to /etc/fstab in operating system.
- ◆ See the *Installation and Configuration Guide* for steps for hybrid disk configuration.

Adjusting OS startup time

It is necessary to configure the time from power-on of each node in the cluster to the server operating system startup to be longer than the following:

- ◆ The time from power-on of the shared disks to the point they become available.
- ◆ Heartbeat timeout time

See the *Installation and Configuration Guide* for configuration steps.

Verifying the network settings

- ◆ The network used by Interconnect or Mirror disk connect is checked. It checks by all the servers in a cluster.
- ◆ See the *Installation and Configuration Guide* for configuration steps.

Ipmitool and OpenIPMI

- ◆ The following functions use ipmitool or OpenIPMI.
 - Final Action at Activation Failure / Deactivation Failure
 - Monitor resource action upon failure
 - User space monitor
 - Shutdown stall monitor
 - Forced Stop
 - Chassis ID lamp
- ◆ ipmitool and OpenIPMI do not come with ExpressCluster. You need to download and install the rpm packages for ipmitool and OpenIPMI.
- ◆ Users are responsible for making decisions and assuming responsibilities. NEC does not support or assume any responsibilities for:
 - Inquires about ipmitool and OpenIPMI themselves.
 - Tested operation of ipmitool and OpenIPMI
 - Malfunction of ipmitool and OpenIPMI or error caused by such malfunction.
 - Inquiries about whether or not ipmitool and OpenIPMI are supported by servers.

- ◆ Check whether or not your server (hardware) supports ipmiutil and OpenIPMI in advance.
- ◆ Note that even if the machine complies with ipmi standard as hardware, ipmiutil and OpenIPMI may not run if you actually try to run them.
- ◆ If you are using a software program for server monitoring provided by a server vendor, do not choose ipmi as a monitoring method for user space monitor resource and shutdown stall monitor. Because these software programs for server monitoring and ipmiutil both use BMC (Baseboard Management Controller) on the server, a conflict occurs preventing successful monitoring.

User mode monitor resource (monitoring method: softdog)

- ◆ When softdog is selected as a monitoring method, make sure to set heartbeat that comes with OS not to start.
- ◆ When it sets softdog in a monitor method in SUSE LINUX 10/11, it is impossible to use with an i8xx_tco driver. When an i8xx_tco driver is unnecessary, please make it the setting that i8xx_tco is not loaded.

Log collection

- ◆ The designated function of the generation of the syslog does not work by a log collection function in SUSE LINUX 10/11. The reason is because the suffixes of the syslog are different.
Please change setting of rotate of the syslog as follows to use the appointment of the generation of the syslog of the log collection function.
- ◆ Please comment out “compress” and “date ext” of the /etc/logrotate.d/syslog file.

nsupdate and nslookup

- ◆ The following functions use nsupdate and nslookup.
 - Dynamic DNS resource of group resource (ddns)
 - Dynamic DNS monitor resource of monitor resource (ddnsw)
- ◆ ExpressCluster does not include nsupdate and nslookup. Therefore, install the rpm files of nsupdate and nslookup, in addition to the ExpressCluster installation.
- ◆ NEC does not support the items below regarding nsupdate and nslookup. Use nsupdate and nslookup at your own risk.
 - Inquiries about nsupdate and nslookup
 - Guaranteed operations of nsupdate and nslookup
 - Malfunction of nsupdate or nslookup or failure caused by such a malfunction
 - Inquiries about support of nsupdate and nslookup on each server

Notes when creating ExpressCluster configuration data

Notes when creating a cluster configuration data and before configuring a cluster system is described in this section.

Environment variable

The following scripts cannot be executed under the environment where more than 255 environmental variables are set. When using the following function of resource, set the number of environmental variables less than 256.

Start/Stop script executed by EXEC resource when activating/deactivating

Script executed by Custom monitor Resource when monitoring

Script before final action after the group resource or the monitor resource error is detected.

Force stop function, chassis identify lamp linkage

When using forced stop function or chassis identify lamp linkage, settings of BMC IP address, user name and password of each server are necessary. Use definitely the username to which the password is set.

Server reset, server panic and power off

When ExpressCluster performs “Server Reset”, “Server Panic,” or “Server power off”, servers are not shut down normally. Therefore, the following may occur.

- ◆ Damage to a mounted file system
- ◆ Lost of unsaved data
- ◆ Suspension of OS dump collection

“Server reset” or “Server panic” occurs in the following settings:

Action at an error occurred when activating/inactivating group resources

- Sysrq Panic
- Keepalive Reset
- Keepalive Panic
- BMC Reset
- BMC Power Off
- BMC Power Cycle
- BMC NMI

- ◆ Final action at detection of an error in monitor resource
 - Sysrq Panic
 - Keepalive Reset
 - Keepalive Panic
 - BMC Reset
 - BMC Power Off
 - BMC Power Cycle
 - BMC NMI
- ◆ Action at detection of user space monitor timeout
 - Monitoring method softdog

- Monitoring method ipmi
- Monitoring method keepalive

Note: “Server panic” can be set only when the monitoring method is “keepalive.”

- ◆ Shutdown stall mentoring
 - Monitoring method softdog
 - Monitoring method ipmi
 - Monitoring method keepalive

Note: “Server panic” can be set only when the monitoring method is “keepalive.”

- ◆ Operation of Forced Stop
 - BMC reset
 - BMC power off
 - BMC cycle
 - BMC NMI

Final action for group resource deactivation error

If you select **No Operation** as the final action when a deactivation error is detected, the group does not stop but remains in the deactivation error status. Make sure not to set **No Operation** in the production environment.

Stack size of the application executed by EXEC resource

Exec resource is executed while the stack size is configured as 2MB. Thus, if an application which is started from exec resource requires the stack size of more than 2MB, stack overflow occurs.

If stack overflow occurs, configure the stack size before starting the application.

1. If you select **Script created with this product**
Please change stack size using ulimit command before the application is executed.
2. If you select **User Application** (Do not use this mode)
Please select Script created with this product and edit script file to execute the application by the script. Also, please change stack size using ulimit command before the application is executed.

Example of start script (start.sh)

```
-----
#!/bin/sh
#####
#*                start.sh                *
#####
```

```
ulimit -s unlimited # Change stack size (unlimited)
```

```
" the application to be executed"
```

When you will change scripts for exec resource, please refer to Reference Guide Section II “Chapter 5 Group resource details – Understanding EXEC resources”.

Verifying raw device for VxVM

Section II Installing ExpressCluster

Check the raw device of the volume raw device in advance:

1. Import all disk groups which can be activated on one server and activate all volumes before installing ExpressCluster.
2. Run the command below:

```
# raw -qa
```

```
/dev/raw/raw2: bound to major 199, minor 2
```

```
/dev/raw/raw3: bound to major 199, minor 3
```

(A)

(B)

Example: Assuming the disk group name and volume name are:

- Disk group name: dg1
- Volume name under dg1: vol1, vol2

3. Run the command below:

```
# ls -l /dev/vx/dsk/dg1/
```

```
brw----- 1 root root 199, 2 May 15 22:13 vol1
```

```
brw----- 1 root root 199, 3 May 15 22:13 vol2
```

(C)

4. Confirm that major and minor numbers are identical between (B) and (C).

Never use these raw devices (A) as an ExpressCluster disk heartbeat resource, raw resource, raw monitor resource, disk resource for which the disk type is not VxVM, or disk monitor resource for which the monitor method is not READ(VxVM).

Selecting mirror disk file system

Following is the currently supported file systems:

- ◆ ext3
- ◆ xfs
- ◆ reiserfs
- ◆ jfs
- ◆ vxfs

Selecting hybrid disk file system

The following are the currently supported file systems:

ext3

reiserfs

RAW monitoring of disk monitor resources

- ◆ When raw monitoring of disk monitor resources is set up, partitions cannot be monitored if they have been or will possibly be mounted. These partitions cannot be monitored even if you set device name to “whole device” (device indicating the entire disks).
- ◆ Allocate a partition dedicated to monitoring and set up the partition to use the raw monitoring of disk monitor resources.

Delay warning rate

If the delay warning rate is set to 0 or 100, the following can be achieved:

- ◆ When 0 is set to the delay monitoring rate

An alert for the delay warning is issued at every monitoring.

By using this feature, you can calculate the polling time for the monitor resource at the time the server is heavily loaded, which will allow you to determine the time for monitoring time-out of a monitor resource.

- ◆ When 100 is set to the delay monitoring rate

The delay warning will not be issued.

Be sure not to set a low value, such as 0%, except for a test operation.

Disk monitor resource (monitoring method TUR)

- ◆ You cannot use the TUR methods on a disk or disk interface (HBA) that does not support the Test Unit Ready (TUR) and SG_IO commands of SCSI. Even if your hardware supports these commands, consult the driver specifications because the driver may not support them.
- ◆ S-ATA disk interface may be recognized as IDE disk interface (hd) or SCSI disk interface (sd) by OS depending on disk controller type and distribution. When it is recognized as IDE interface, all TUR methods cannot be used. If it is recognized as SCSI disk interface, TUR (legacy) can be used. Note that TUR (generic) cannot be used.
- ◆ TUR methods burdens OS and disk load less compared to Read methods.
- ◆ In some cases, TUR methods may not be able to detect errors in I/O to the actual media.

WebManager reload interval

- ◆ Do not set the “Reload Interval” in the WebManager tab for less than 30 seconds.

LAN heartbeat settings

- ◆ You need to set at least one LAN heartbeat resource. It is recommended to set two or more LAN heartbeat resources.
- ◆ It is recommended to set both LAN heartbeat resource and kernel mode LAN heartbeat resource together.

Kernel mode LAN heartbeat resource settings

- ◆ It is recommended to use both LAN heartbeat resource and kernel mode LAN heartbeat resource for distribution kernel of which kernel mode LAN heartbeat can be used.
- ◆ It is recommended to register interconnect-dedicated LAN and public LAN as LAN heartbeat resource and kernel mode LAN heartbeat resource. (Registering more than two LAN heartbeat resources and kernel mode LAN heartbeat resources is recommended.)

COM heartbeat resource settings

- ◆ It is recommended to use a COM heartbeat resource if your environments allows. This is because using COM heartbeat resource prevents activating both systems when the network is disconnected.

Double-byte character set that can be used in script comments

- ◆ Scripts edited in Linux environment are dealt as EUC code, and scripts edited in Windows environment are dealt as Shift-JIS code. In case that other character codes are used, character corruption may occur depending on environment.

Failover exclusive attribute of virtual machine group

- ◆ When setting virtual machine group, do not set **Normal** or **Absolute** to **Failover exclusive attribute**.

After start operating ExpressCluster

Notes on situations you may encounter after start operating ExpressCluster are described in this section.

Error message in the load of the mirror driver in the udev environment

In the load of the mirror driver in the udev environment, logs like the following may be recorded into the message file:

```
kernel: [I] <type: liscal><event: 141> NMP1 device does not exist.
(liscal_make_request)
kernel: [I] <type: liscal><event: 141> - This message can be recorded
on udev environment when liscal is initializing NMPx.
kernel: [I] <type: liscal><event: 141> - Ignore this and following
messages 'Buffer I/O error on device NMPx' on udev environment.
kernel: Buffer I/O error on device NMP1, logical block 0
```

```
kernel: <liscal liscal_make_request> NMP1 device does not exist.
kernel: Buffer I/O error on device NMP1, logical block 112
```

This phenomenon is not abnormal.

When you want to evade the output of the error message, please add the following preference to `/etc/udev/rules.d/` subordinates.

filename: 50-liscal-udev.rules

```
ACTION=="add", DEVPATH=="block/NMP*", OPTIONS+="ignore_device"
```

File operating utility on X-Window

Some of the file operating utilities (copying and moving files and directories via GUI) on X-Window perform the following:

- ◆ Checks if the block device is usable.
- ◆ Mounts the file system if there is any that can be mounted.

Make sure not to use file operating utility that perform above operations. They may cause problem to the operation of ExpressCluster.

Messages displayed when loading a driver

When loading a mirror driver, messages like the following may be displayed at the console and/or syslog. However, this is not an error.

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

(Any character strings are set to *xxxxx*.)

And also, when loading the *clpka* or *clpkhb* driver, messages like the following may be displayed on the console and/or *syslog*. However, this is not an error.

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

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

(Any character strings are input into *xxxxx*.)

Messages written to *syslog* when multiple mirror disk resources or hybrid disk resources are used

When more than two mirror disk resources or hybrid disk resources are configured on a cluster, the following messages may be written to the OS message files when the resources are activated.

This occurs by a *fsck* command function (function to access a device block which is not a target of *fsck*) on some distributions.

```
kernel: [I] <type: liscal><event: 144> NMPx I/O port has been closed,
mount(0), io(0).
kernel: [I] <type: liscal><event: 144> - This message can be recorded
by fsck command when NMPx becomes active.
kernel: [I] <type: liscal><event: 144> - This message can be recorded
on hotplug service starting when NMPx is not active.
kernel: [I] <type: liscal><event: 144> - Ignore this and following
messages 'Buffer I/O error on device NMPx' on such environment.
kernel: Buffer I/O error on device /dev/NMPx, logical block xxxx
```

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

This is not a problem for ExpressCluster. If this causes any problem such as heavy use of message files, change the following settings of mirror resources or hybrid disk resources.

- Select "Not Execute" on "fsck action before mount"
- Select "Execute" on "fsck Action When Mount Failed"

IPMI message

When you are using ipmi for user mode monitor resources, the following kernel module warning log is recorded many times in the syslog.

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

When you want to prevent this log from being recorded, rename /dev/ipmikcs.

Limitations during the recovery operation

Do not control the following commands, clusters and groups by the WebManager while recovery processing is changing (reactivation → failover → last operation), if a group resource is specified as a recovery target and when a monitor resource detects an error.

- ◆ Stop and suspend of a cluster
- ◆ Start, stop, moving of a group

If these operations are controlled at the transition to recovering due to an error detected by a monitor resource, the other group resources in the group may not be stopped.

Even if a monitor resource detects an error, it is possible to control the operations above after the last operation is performed.

Executable format file and script file not described in manuals

Executable format files and script files which are not described in Chapter 4, "ExpressCluster command reference" in the *Reference Guide* exist under the installation directory. Do not run these files on any system other than ExpressCluster. The consequences of running these files will not be supported.

Message of kernel page allocation error

When using the Replicator on the TurboLinux 10 Server, the following message may be recorded in syslog. However, it may not be recorded depending on the physical memory size and I/O load.

```
kernel: [kernel Module Name]: page allocation failure. order:X,  
mode:0xXX
```

When this message is recorded, you need to change the kernel parameter described below. By using the sysctl command, make the settings to change the parameter when starting OS.

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

The maximum value that can be set to min_free_kbyte is different depending on the physical memory size installed on the server. Make the settings by referring to the table below:

Physical memory size (Mbyte)	Maximum value (Mbyte)
1024	1024

2048	1448
4096	2048
8192	2896
16384	4096

Cache swell by a massive I/O to mirror disk resource and hybrid disk resource

- ◆ In case that a massive amount of write over the disk capability to the mirror disk resource or the hybrid disk resource are executed, even though the mirror connection is alive, the control from write may not return or memory allocation failure may occur.

In case that a massive amount of I/O requests over transaction performance exist, and then the file system ensure a massive amount of cache and the cache or the memory for the user space (HIGHMEM zone) are insufficient, the memory for the kernel space (NORMAL zone) may be used.

Change the settings so that the parameter will be changed at OS startup by using `sysctl` or other commands.

```
/proc/sys/vm/lower_zone_protection
```

- ◆ In case that a massive amount of accesses to the mirror disk resource or the hybrid disk resource are executed, it may take much time before the cache of the file systems is written out to the disks when unmounting at disk resource deactivation.

At times like this, set the timeout interval of unmount longer enough so that the writing to the disks will successfully complete.

Messages when collecting logs

When collecting logs, the message described below is displayed at the console, but this is not an error. Logs are collected successfully.

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

(“hd#” is replaced with the device name of IDE.)

```
kernel: Warning: /proc/ide/hd?/settings interface is obsolete, and
will be removed soon!
```

Cluster shutdown and reboot

When using a mirror disk resource or a hybrid disk resource, do not execute cluster shutdown or cluster shutdown reboot from the `clpstdn` command or the WebManager while a group is being activated.

A group cannot be deactivated while a group is being activated. Therefore, OS may be shut down in the state that mirror disk resource or hybrid disk resources is not deactivated successfully and a mirror break may occur.

Shutdown and reboot of individual server

When using a mirror disk and a hybrid disk resource, do not shut down the server or run the shutdown reboot command from the `clpdown` command or the WebManager while activating the group.

A group cannot be deactivated while a group is being activated. Therefore, OS may be shut down and a mirror break may occur in the state that mirror disk resources and hybrid disk resources are not deactivated successfully.

Scripts for starting/stopping ExpressCluster services

Errors occur in starting/stopping scripts as follows:

- ◆ After installing ExpressCluster (For SUSE Linux)
When a server shutdown, the error occurs in the following stopping scripts. There is no problem for the error because services have not started.
 - `clusterpro_alertsync`
 - `clusterpro_webmgr`
 - `clusterpro`
 - `clusterpro_md`
 - `clusterpro_trn`
 - `clusterpro_evt`
- ◆ Before start operationg ExpressCluster
When a server start up, the error occurs in the following starting scripts. There is no problem for the error because cluster configuration data has not uploaded.
 - `clusterpro_md`
- ◆ After start operating ExpressCluster (For SUSE Linux)
When mirror disk resources and hybrid disk resources are not used, the error occurs in stopping scripts at OS shutdown. There is no problem for the error because mirror agent has not started.
 - `clusterpro_md`
- ◆ OS shutdown after stopping services manually (Fro SUSE Linux)
After stopping services manually, the error occurs in the following stopping scripts at OS shutdown. There is no problem for the error because services have already stopped.
 - `clusterpro`
 - `clusterpro_md`

At following case, the script to terminate ExpressCluster services may be executed in the wrong order.

- ◆ ExpressCluster services may be terminated in the wrong order at OS shutdown if all of ExpressCluster services are disabled. This problem is caused by failure in termination process for the service has been already disabled.
As long as the system shutdown is executed by WebManger or clpstdn command, there is no problem even if the services is terminated in the wrong order. But, any other problem may not be happened by wrong order termination.

Service startup time

ExpressCluster services might take a while to start up, depending on the wait processing at startup.

- ◆ clusterpro_evt
Servers other than the master server wait up to two minutes for configuration data to be downloaded from the master server. Downloading usually finishes within several seconds if the master server is already operating. The master server does not have this wait process.
- ◆ clusterpro_trn
There is no wait process. This process usually finishes within several seconds.
- ◆ clusterpro_md
This service starts up only when the mirror or hybrid disk resources exist. The system waits up to one minute for the mirror agent to normally start up. This process usually finishes within several seconds.
- ◆ clusterpro
Although there is no wait process, ExpressCluster might take several tens of seconds to start up. This process usually finishes within several seconds.
- ◆ clusterpro_webmgr
There is no wait process. This process usually finishes within several seconds.
- ◆ clusterpro_alertsync
There is no wait process. This process usually finishes within several seconds.

In addition, the system waits for cluster activation synchronization after the ExpressCluster daemon is started. By default, this wait time is five minutes.

For details, see Chapter 10, “The system maintenance information” in the *Reference Guide*.

Scripts in EXEC resources

EXEC resource scripts of group resources stored in the following location.

```
/opt/nec/clusterpro/scripts/group-name/resource-name/
```

The following cases, old EXEC resource scripts are not deleted automatically.

- When the EXEC resource is deleted or renamed
- When a group that belongs to the EXEC resource is deleted or renamed

Old EXEC resource scripts can be deleted when unnecessary.

Monitor resources that monitoring timing is “Active”

When monitor resources that monitoring timing is “Active” have suspended and resumed, the following restriction apply:

- ◆ In case stopping target resource after suspending monitor resource, monitor resource becomes suspended. As a result, monitoring restart cannot be executed.
- ◆ In case stopping or starting target resource after suspending monitor resource, monitoring by monitor resource starts when target resource starts.

Notes on the WebManager

- ◆ The information displayed on the WebManager does not necessarily show the latest status. If you want to get the latest information, click the **Reload** button.
- ◆ If the problems such as server shutdown occur while the WebManager is getting the information, acquiring information may fail and a part of object may not be displayed correctly. Wait for the next automatic update or click the **Reload** button to reacquire the latest information.
- ◆ When using a browser on Linux, a dialog box may be displayed behind the window managers depending on the combination of the managers. Change the window by pressing the **ALT + TAB** keys.
- ◆ Collecting logs of ExpressCluster cannot be executed from two or more WebManager simultaneously.
- ◆ If the WebManager is operated in the state that it cannot communicate with the connection destination, it may take a while until the control returns.
- ◆ If you move the cursor out of the browser in the state that the mouse pointer is displayed as a wristwatch or hourglass, the cursor may be back to an arrow.
- ◆ When going through the proxy server, make the settings for the proxy server be able to relay the port number of the WebManager.
- ◆ When updating ExpressCluster, close the browser. Clear the Java cache and open the browser.

Notes on the Builder (Config mode of Cluster Manager)

- ◆ ExpressCluster does not have the compatibility of the cluster configuration data with the following products.
 - Builder for Linux other than ExpressCluster X 3.0 for Linux
- ◆ Closing the Web browser (by clicking **Exit** from the menu) discards the edited data. Even if the configuration is changed, the dialog box to confirm to save is not displayed. When you need to save the edited data, select **File** from the menu of the Builder and click **Export** before terminating.
- ◆ Reloading the Web browser (by selecting **Refresh** button from the menu or tool bar) discards the current editing data. Even if the configuration is changed, the dialog box to confirm to save is not displayed. When you need to save the editing data, select **File** from the menu bar of the Builder and click **Export** before reloading.
- ◆ When creating the cluster configuration data using the Builder, do not enter the value starting with 0 on the text box. For example, if you want to set 10 seconds for a timeout value, enter “10” but not “010.”

Changing the partition size of mirror disks and hybrid disk resources

When changing the size of mirror partitions after the operation is started, see “Changing offset or size of a partition on mirror resource” in Chapter 10 “The system maintenance information” in the *Reference Guide*.

Chapter 6 **Upgrading ExpressCluster**

This chapter provides information on how to upgrade ExpressCluster.
This chapter covers:

- How to update from ExpressCluster X 2.0 or 2.1..... 104

How to update from ExpressCluster X 2.0 or 2.1

How to upgrade from X2.0 or X2.1

Install the ExpressCluster Server RPM as root user.

1. Disable the services by running the **chkconfig --del name** in the following order on all the servers. Specify one of the following services in *name*.

clusterpro_alertsync

clusterpro_webmgr

clusterpro

clusterpro_md

clusterpro_trn

clusterpro_evt

2. Shut down and reboot the cluster by using WebManager or the `clpstdn` command.
3. Mount the installation CD-ROM media.
4. Confirm that ExpressCluster services are not running, and then install the package file by executing the `rpm` command. The RPM for installation is different depending on architecture.

In the CD-ROM, move to `/Linux/3.0/en/server` and run the following:

```
rpm -U expresscls-<version>.<architecture>.rpm
```

For architecture, there are `i686`, `x86_64` and `ppc64`. Select architecture according to the system requirements of the machine where ExpressCluster is installed. Architecture can be verified by the `arch` command.

ExpressCluster is installed in the following directory. Note that if you change this directory you cannot uninstall ExpressCluster.

Installation directory: `/opt/nec/clusterpro`

5. After completing installation, unmount the installation CD-ROM media, and remove it.
6. Enable the services by running the **chkconfig --add name** in the following order. Specify one of the following services in *name*. For SUSE Linux, run the command with the **-force** option.

clusterpro_evt

clusterpro_trn

clusterpro_webmgr

clusterpro_alertsync
7. Repeat the steps 3-6 on all the servers.
8. Reboot all the servers that constitute the cluster.
9. Register the license. For details on registering license, see “Chapter 4 Registering the license” in the *Installation and Configuration Guide*.
10. Connect the WebManager to one of the servers of the cluster.
11. Start the Builder from the connected WebManager. For details on how to start the online Builder, see the *Installation and Configuration Guide*.

12. Confirm that all servers of the cluster are started, and then upload the configuration data from the online Builder. For details on how to operate the online Builder, see the *Reference Guide*.
13. Enable the services in the following order by running the `chkconfig --add name` command. Specify the following services on *name*.
clusterpro_md
clusterpro
14. Perform step 14 on all the servers.
15. Run **Restart Manager** on the WebManager.
16. Run **Start Mirror Agent** on the WebManager.
17. Restart the browser connecting the WebManager.
18. Run **Start Cluster** on the WebManager.

Appendix

- Appendix A Glossary
- Appendix B Index

Appendix A. Glossary

Cluster partition	A partition on a mirror disk. Used for managing mirror disks. (Related term: Disk heartbeat partition)
Interconnect	A dedicated communication path for server-to-server communication in a cluster. (Related terms: Private LAN, Public LAN)
Virtual IP address	IP address used to configure a remote cluster.
Management client	Any machine that uses the WebManager to access and manage a cluster system.
Startup attribute	A failover group attribute that determines whether a failover group should be started up automatically or manually when a cluster is started.
Shared disk	A disk that multiple servers can access.
Shared disk type cluster	A cluster system that uses one or more shared disks.
Switchable partition	A disk partition connected to multiple computers and is switchable among computers. (Related terms: Disk heartbeat partition)
Cluster system	Multiple computers are connected via a LAN (or other network) and behave as if it were a single system.
Cluster shutdown	To shut down an entire cluster system (all servers that configure a cluster system).
Active server	A server that is running for an application set. (Related term: Standby server)
Secondary server	A destination server where a failover group fails over to during normal operations. (Related term: Primary server)
Standby server	A server that is not an active server. (Related term: Active server)
Disk heartbeat partition	A partition used for heartbeat communication in a shared disk type cluster.
Data partition	A local disk that can be used as a shared disk for switchable partition. Data partition for mirror disks and hybrid disks. (Related term: Cluster partition)
Network partition	All heartbeat is lost and the network between servers is partitioned. (Related terms: Interconnect, Heartbeat)

Node	A server that is part of a cluster in a cluster system. In networking terminology, it refers to devices, including computers and routers, that can transmit, receive, or process signals.
Heartbeat	Signals that servers in a cluster send to each other to detect a failure in a cluster. (Related terms: Interconnect, Network partition)
Public LAN	A communication channel between clients and servers. (Related terms: Interconnect, Private LAN)
Failover	The process of a standby server taking over the group of resources that the active server previously was handling due to error detection.
Failback	A process of returning an application back to an active server after an application fails over to another server.
Failover group	A group of cluster resources and attributes required to execute an application.
Moving failover group	Moving an application from an active server to a standby server by a user.
Failover policy	A priority list of servers that a group can fail over to.
Private LAN	LAN in which only servers configured in a clustered system are connected. (Related terms: Interconnect, Public LAN)
Primary (server)	A server that is the main server for a failover group. (Related term: Secondary server)
Floating IP address	Clients can transparently switch one server from another when a failover occurs. Any unassigned IP address that has the same network address that a cluster server belongs to can be used as a floating address.
Master server	The server displayed on top of the Master Server in Cluster Properties in the Builder.
Mirror disk connect	LAN used for data mirroring in mirror disks and hybrid disks. Mirror connect can be used with primary interconnect.
Mirror disk type cluster	A cluster system that does not use a shared disk. Local disks of the servers are mirrored.

Appendix B. Index

A

application monitoring, 33
Applications supported, 58

B

browsers, 62, 63, 64
Builder, 62, 72, 100

C

clock synchronization, 85
cluster object, 42
Cluster shutdown and reboot, 98
cluster system, 16
COM heartbeat, 89
COM heartbeat resource, 92
communication port number, 83
Config mode of Cluster Manager, 100
Corrected information, 69

D

data consistency, 78
delay warning rate, 29
dependent driver, 82
dependent library, 81
detectable and non-detectable errors, 33, 34
disk interfaces, 52
disk size, 64
distribution, 54

E

Enhanced functions, 67
Environment variable, 88
error detection, 15, 20
executable format file, 96
ExpressCluster, 29, 30

F

failover, 23, 29, 34
failover resources, 35
failure monitoring, 27
File operating utility, 94
file system, 78, 90
final action, 89
Force stop function, chassis identify lamp linkage, 88

G

group resource, 89
group resources, 43

H

hardware, 52

hardware configuration, 39, 40, 41
hardware requirements for hybrid disk, 75
hardware requirements for mirror disk, 72
hardware requirements for shared disk, 73
heartbeat resources, 43
High Availability (HA) cluster, 16
How an error is detected, 32
hybrid disk, 81, 86, 100

I

inheriting applications, 22
inheriting cluster resources, 21
inheriting data, 21
integrated WebManager, 64
internal monitoring, 33
IPMI message, 96

J

Java runtime environment, 62, 63, 64

K

kernel, 54
Kernel mode LAN heartbeat and keepalive drivers, 82
kernel mode LAN heartbeat resource, 92

L

LAN heartbeat, 92
log collection, 87

M

memory and disk size, 61, 62, 63
memory size, 64
message of kernel page allocation error, 96
messages displayed when loading a driver, 94
messages when collecting logs, 97
mirror disk, 78, 85
mirror driver, 82
modules, 30
monitor resources, 44
monitor resources that monitoring timing is, 99
monitored and non-monitored errors, 33

N

network interfaces, 53
Network partition, 21
Network partition resolution resources, 43
network settings, 86
NIC device name, 85
NIC link up/down monitor resource, 76

O

O_DIRECT, 78
operating systems, 62, 63, 64

OS startup time, 86

R

raw device, 89
raw monitor resources, 90
RAW monitor resources, 82
reload interval, 92
resource, 29, 43

S

script file, 96
scripts for starting/stopping ExpressCluster services, 98
scripts in EXEC resources, 99
server monitoring, 32
server requirements, 52
server reset and panic, 88
shared disk, 85
shutdown and reboot of individual server, 98

single point of failure, 24
software, 54
software configuration, 29, 30
supported operating systems, 72
system configuration, 36

T

TUR, 92

U

user mode monitor resource, 87

W

WebManager, 63, 64, 72, 100
write function, 77