
Pacemaker Explained

Release 3.0.0

the Pacemaker project contributors

Jan 09, 2025

CONTENTS

1	Abstract	3
2	Table of Contents	5
2.1	Introduction	5
2.1.1	The Scope of this Document	5
2.1.2	What Is Pacemaker?	5
2.2	Host-Local Configuration	10
2.2.1	Configuration Value Types	10
2.2.2	Local Options	12
2.3	Cluster-Wide Configuration	19
2.3.1	Configuration Layout	19
2.3.2	Option Precedence	20
2.3.3	CIB Properties	21
2.3.4	Cluster Options	22
2.4	Nodes	29
2.4.1	Cluster nodes	29
2.4.2	Pacemaker Remote nodes	30
2.4.3	Defining a Node	30
2.4.4	Quorum-only Nodes	31
2.4.5	Node Attributes	31
2.4.6	Tracking Node Health	34
2.5	Resources	36
2.5.1	Resource Standards	36
2.5.2	Resource Properties	38
2.5.3	Resource Options	39
2.5.4	Pacemaker Remote Resources	44
2.6	Resource Operations	45
2.6.1	Operation Properties	46
2.6.2	Monitoring Resources for Failure	49
2.6.3	Custom Recurring Operations	50
2.6.4	Setting Global Defaults for Operations	50
2.6.5	When Implicit Operations Take a Long Time	50
2.6.6	Multiple Monitor Operations	50
2.6.7	Disabling a Monitor Operation	51
2.6.8	Handling Resource Failure	52
2.6.9	Reloading an Agent After a Definition Change	53
2.6.10	Migrating Resources	54
2.7	Resource Constraints	55
2.7.1	Deciding Which Nodes a Resource Can Run On	55
2.7.2	Specifying the Order in which Resources Should Start/Stop	58

2.7.3	Placing Resources Relative to other Resources	60
2.7.4	Resource Sets	62
2.7.5	Ordering Sets of Resources	63
2.7.6	Colocating Sets of Resources	67
2.7.7	External Resource Dependencies	69
2.8	Fencing	70
2.8.1	What Is Fencing?	70
2.8.2	Why Is Fencing Necessary?	70
2.8.3	Fence Devices	70
2.8.4	Fence Agents	71
2.8.5	When a Fence Device Can Be Used	71
2.8.6	Limitations of Fencing Resources	71
2.8.7	Special Meta-Attributes for Fencing Resources	72
2.8.8	Special Instance Attributes for Fencing Resources	72
2.8.9	Default Check Type	75
2.8.10	Unfencing	76
2.8.11	Fencing and Quorum	76
2.8.12	Fencing Timeouts	76
2.8.13	Fence Devices Dependent on Other Resources	77
2.8.14	Configuring Fencing	77
2.8.15	Fencing Topologies	84
2.8.16	Remapping Reboots	87
2.9	Collective Resources	87
2.9.1	Groups - A Syntactic Shortcut	87
2.9.2	Clones - Resources That Can Have Multiple Active Instances	89
2.9.3	Bundles - Containerized Resources	95
2.10	Utilization and Placement Strategy	101
2.10.1	Utilization attributes	101
2.10.2	Placement Strategy	102
2.10.3	How Multiple Capacities Combine	103
2.10.4	Order of Resource Assignment	103
2.10.5	Limitations	103
2.11	Rules	104
2.11.1	Rule Options	104
2.11.2	Rule Conditions and Contexts	105
2.11.3	Date/Time Expressions	105
2.11.4	Node Attribute Expressions	109
2.11.5	Resource Type Expressions	111
2.11.6	Operation Type Expressions	111
2.11.7	Using Rules to Determine Resource Location	112
2.11.8	Using Rules to Define Options	115
2.12	Access Control Lists (ACLs)	118
2.12.1	ACL Prerequisites	118
2.12.2	ACL Configuration	118
2.12.3	ACL Roles	118
2.12.4	ACL Targets and Groups	119
2.12.5	ACLs and Pacemaker Remote Nodes	120
2.12.6	ACL Examples	120
2.12.7	ACL Limitations	123
2.13	Alerts	123
2.13.1	Alert Agents	123
2.13.2	Alert Recipients	124
2.13.3	Alert Meta-Attributes	124
2.13.4	Alert Instance Attributes	125

2.13.5	Alert Filters	126
2.14	Reusing Parts of the Configuration	127
2.14.1	Reusing Resource Definitions	127
2.14.2	Reusing Rules, Options and Sets of Operations	131
2.14.3	Tagging Configuration Elements	132
2.15	Status	134
2.15.1	Node State	134
2.15.2	Transient Node Attributes	135
2.15.3	Node History	135
2.16	Multi-Site Clusters and Tickets	138
2.16.1	Challenges for Multi-Site Clusters	138
2.16.2	Conceptual Overview	138
2.16.3	Configuring Ticket Dependencies	140
2.16.4	Managing Multi-Site Clusters	141
2.16.5	For more information	143
2.17	Sample Configurations	143
2.17.1	Empty	143
2.17.2	Simple	143
2.17.3	Advanced Configuration	144
3	Index	147
	Index	149

Configuring Pacemaker Clusters

ABSTRACT

This document definitively explains Pacemaker's features and capabilities, particularly the XML syntax used in Pacemaker's Cluster Information Base (CIB).

TABLE OF CONTENTS

2.1 Introduction

2.1.1 The Scope of this Document

This document is intended to be an exhaustive reference for configuring Pacemaker. To achieve this, it focuses on the XML syntax used to configure the CIB.

For those that are allergic to XML, multiple higher-level front-ends (both command-line and GUI) are available. These tools will not be covered in this document, though the concepts explained here should make the functionality of these tools more easily understood.

Users may be interested in other parts of the [Pacemaker documentation set](#), such as *Clusters from Scratch*, a step-by-step guide to setting up an example cluster, and *Pacemaker Administration*, a guide to maintaining a cluster.

2.1.2 What Is Pacemaker?

Pacemaker is a high-availability *cluster resource manager* – software that runs on a set of hosts (a *cluster of nodes*) in order to preserve integrity and minimize downtime of desired services (*resources*).¹ It is maintained by the [ClusterLabs](#) community.

Pacemaker's key features include:

- Detection of and recovery from node- and service-level failures
- Ability to ensure data integrity by fencing faulty nodes
- Support for one or more nodes per cluster
- Support for multiple resource interface standards (anything that can be scripted can be clustered)
- Support (but no requirement) for shared storage
- Support for practically any redundancy configuration (active/passive, N+1, etc.)
- Automatically replicated configuration that can be updated from any node
- Ability to specify cluster-wide relationships between services, such as ordering, colocation, and anti-colocation
- Support for advanced service types, such as *clones* (services that need to be active on multiple nodes), *promotable clones* (clones that can run in one of two roles), and containerized services

¹ *Cluster* is sometimes used in other contexts to refer to hosts grouped together for other purposes, such as high-performance computing (HPC), but Pacemaker is not intended for those purposes.

- Unified, scriptable cluster management tools

Note: Fencing

Fencing, also known as *STONITH* (an acronym for Shoot The Other Node In The Head), is the ability to ensure that it is not possible for a node to be running a service. This is accomplished via *fence devices* such as intelligent power switches that cut power to the target, or intelligent network switches that cut the target's access to the local network.

Pacemaker represents fence devices as a special class of resource.

A cluster cannot safely recover from certain failure conditions, such as an unresponsive node, without fencing.

Cluster Architecture

At a high level, a cluster can be viewed as having these parts (which together are often referred to as the *cluster stack*):

- **Resources:** These are the reason for the cluster's being – the services that need to be kept highly available.
- **Resource agents:** These are scripts or operating system components that start, stop, and monitor resources, given a set of resource parameters. These provide a uniform interface between Pacemaker and the managed services.
- **Fence agents:** These are scripts that execute node fencing actions, given a target and fence device parameters.
- **Cluster membership layer:** This component provides reliable messaging, membership, and quorum information about the cluster. Currently, Pacemaker supports *Corosync* as this layer.
- **Cluster resource manager:** Pacemaker provides the brain that processes and reacts to events that occur in the cluster. These events may include nodes joining or leaving the cluster; resource events caused by failures, maintenance, or scheduled activities; and other administrative actions. To achieve the desired availability, Pacemaker may start and stop resources and fence nodes.
- **Cluster tools:** These provide an interface for users to interact with the cluster. Various command-line and graphical (GUI) interfaces are available.

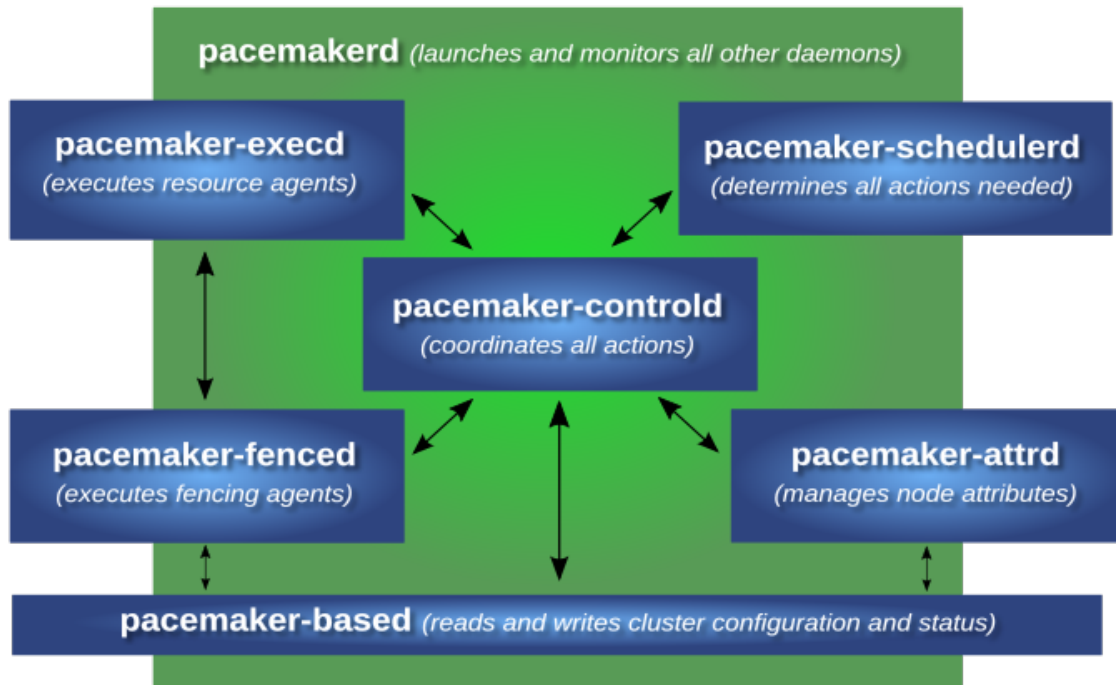
Most managed services are not, themselves, cluster-aware. However, many popular open-source cluster filesystems make use of a common *Distributed Lock Manager* (DLM), which makes direct use of *Corosync* for its messaging and membership capabilities and *Pacemaker* for the ability to fence nodes.

Pacemaker Architecture

Pacemaker itself is composed of multiple daemons that work together:

- `pacemakerd`
- `pacemaker-attd`
- `pacemaker-based`
- `pacemaker-controld`
- `pacemaker-execd`
- `pacemaker-fenced`
- `pacemaker-schedulerd`

Pacemaker internals



ClusterLabs

Pacemaker's main process (`pacemakerd`) spawns all the other daemons, and respawns them if they unexpectedly exit.

The *Cluster Information Base* (CIB) is an XML representation of the cluster's configuration and the state of all nodes and resources. The *CIB manager* (`pacemaker-based`) keeps the CIB synchronized across the cluster, and handles requests to modify it.

The *attribute manager* (`pacemaker-attrd`) maintains a database of attributes for all nodes, keeps it synchronized across the cluster, and handles requests to modify them. These attributes are usually recorded in the CIB.

Given a snapshot of the CIB as input, the *scheduler* (`pacemaker-schedulerd`) determines what actions are necessary to achieve the desired state of the cluster.

The *local executor* (`pacemaker-execd`) handles requests to execute resource agents on the local cluster node, and returns the result.

The *fencer* (`pacemaker-fenced`) handles requests to fence nodes. Given a target node, the fencer decides which cluster node(s) should execute which fencing device(s), and calls the necessary fencing agents (either directly, or via requests to the fencer peers on other nodes), and returns the result.

The *controller* (`pacemaker-controld`) is Pacemaker's coordinator, maintaining a consistent view of the cluster membership and orchestrating all the other components.

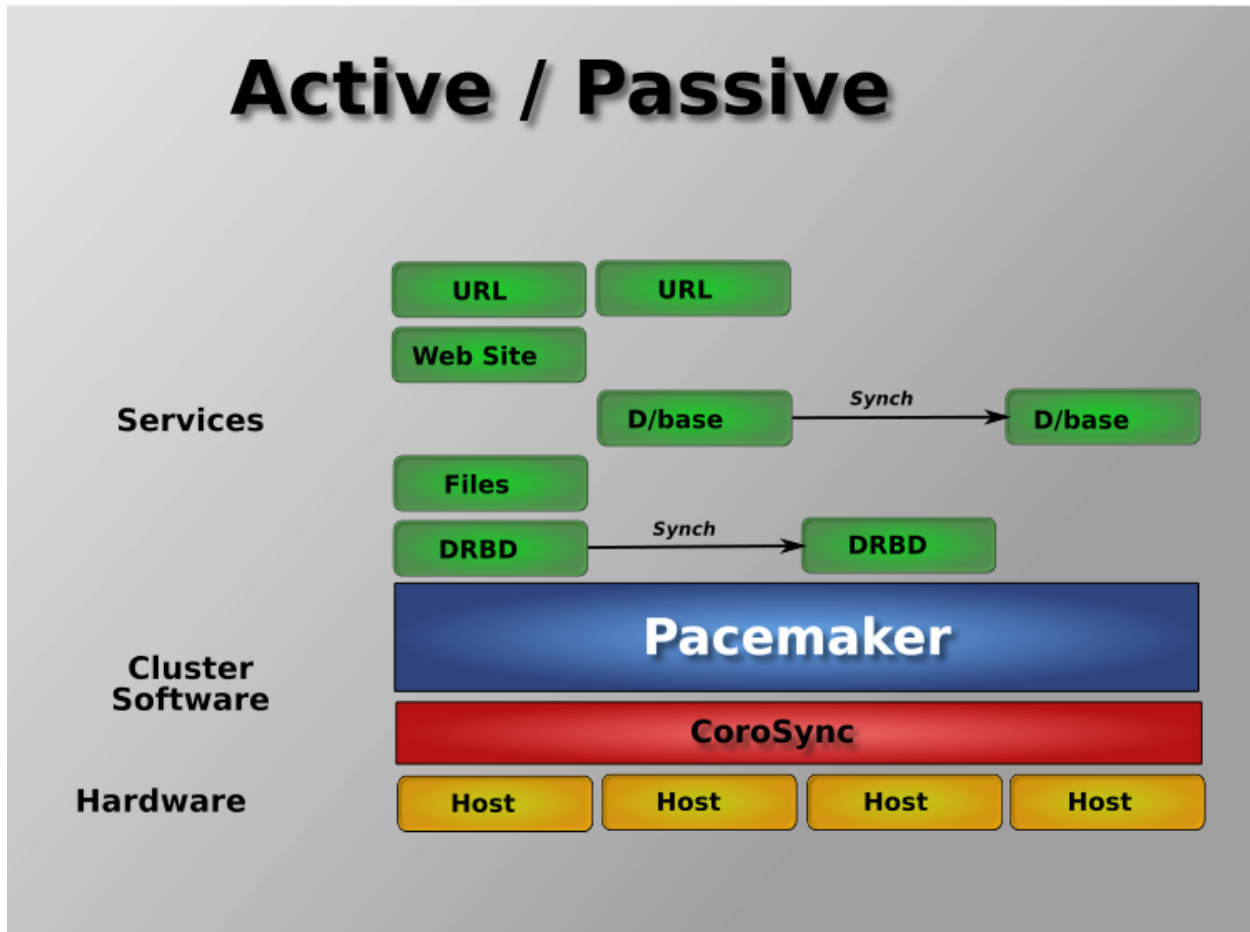
Pacemaker centralizes cluster decision-making by electing one of the controller instances as the *Designated Controller* (DC). Should the elected DC process (or the node it is on) fail, a new one is quickly established. The DC responds to cluster events by taking a current snapshot of the CIB, feeding it to the scheduler, then

asking the executors (either directly on the local node, or via requests to controller peers on other nodes) and the fence to execute any necessary actions.

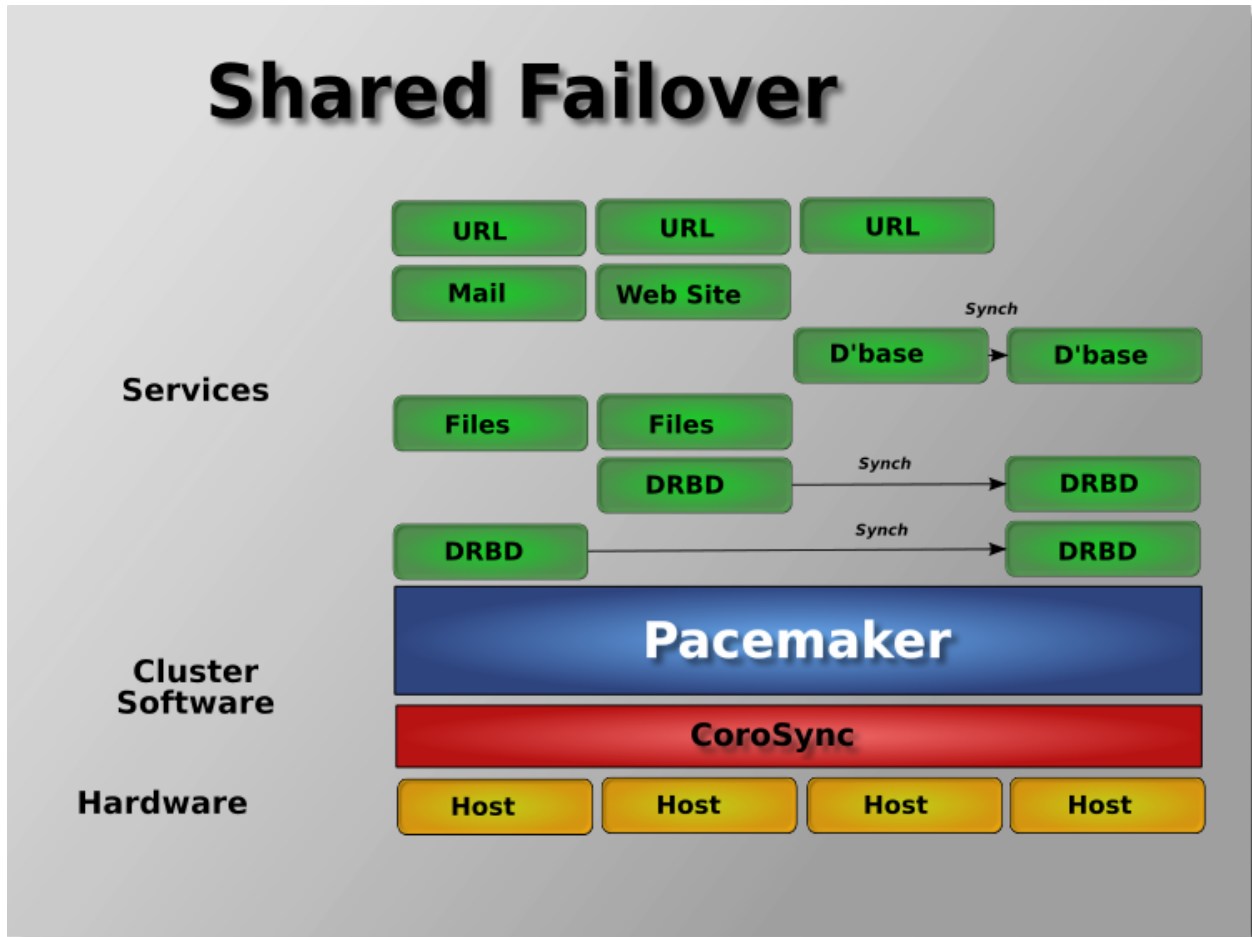
Node Redundancy Designs

Pacemaker supports practically any node redundancy configuration including *Active/Active*, *Active/Passive*, *N+1*, *N+M*, *N-to-1*, and *N-to-N*.

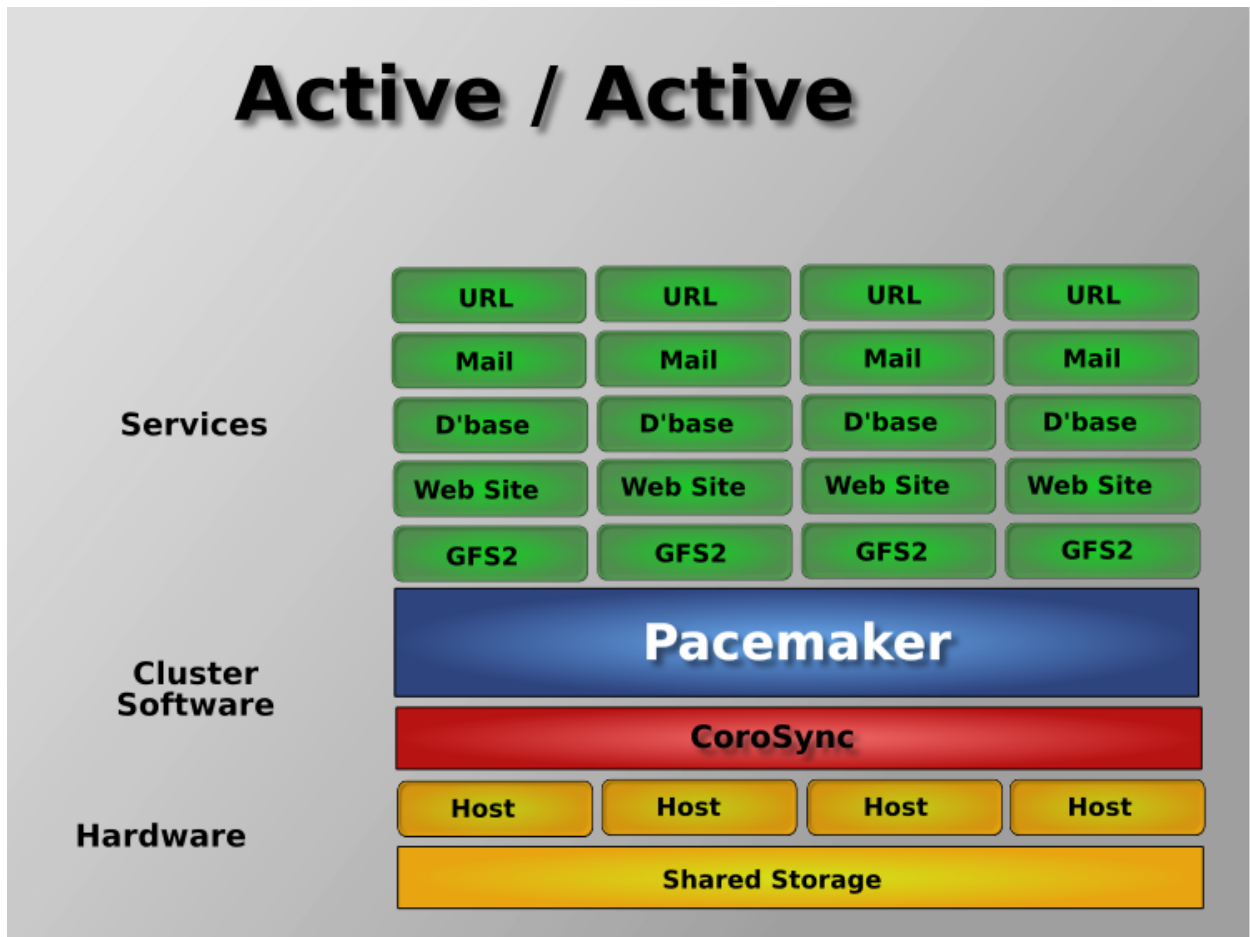
Active/passive clusters with two (or more) nodes using Pacemaker and DRBD are a cost-effective high-availability solution for many situations. One of the nodes provides the desired services, and if it fails, the other node takes over.



Pacemaker also supports multiple nodes in a shared-failover design, reducing hardware costs by allowing several active/passive clusters to be combined and share a common backup node.



When shared storage is available, every node can potentially be used for failover. Pacemaker can even run multiple copies of services to spread out the workload. This is sometimes called N-to-N redundancy.



2.2 Host-Local Configuration

Note: Directory and file paths below may differ on your system depending on your Pacemaker build settings. Check your Pacemaker configuration file to find the correct paths.

2.2.1 Configuration Value Types

Throughout this document, configuration values will be designated as having one of the following types:

Table 1: Configuration Value Types

Type	Description
boolean	Case-insensitive text value where 1, yes, y, on, and true evaluate as true and 0, no, n, off, false, and unset evaluate as false
date/time	Textual timestamp like Sat Dec 21 11:47:45 2013
duration	A nonnegative time duration, specified either like a <i>timeout</i> or an ISO 8601 <i>duration</i> . A duration may be up to approximately 49 days but is intended for much smaller time periods.

Continued on next page

Table 1 – continued from previous page

Type	Description
enumeration	Text that must be one of a set of defined values (which will be listed in the description)
epoch_time	Time as the integer number of seconds since the Unix epoch, 1970-01-01 00:00:00 +0000 (UTC).
id	A text string starting with a letter or underbar, followed by any combination of letters, numbers, dashes, dots, and/or underbars; when used for a property named <code>id</code> , the string must be unique across all <code>id</code> properties in the CIB
integer	32-bit signed integer value (-2,147,483,648 to 2,147,483,647)
ISO 8601	An ISO 8601 date/time.
nonnegative integer	32-bit nonnegative integer value (0 to 2,147,483,647)
percentage	Floating-point number followed by an optional percent sign ('%')
port	Integer TCP port number (0 to 65535)
range	A range may be a single nonnegative integer or a dash-separated range of nonnegative integers. Either the first or last value may be omitted to leave the range open-ended. Examples: 0, 3-, -5, 4-6.
score	A Pacemaker score can be an integer between -1,000,000 and 1,000,000, or a string alias: <code>INFINITY</code> or <code>+INFINITY</code> is equivalent to 1,000,000, <code>-INFINITY</code> is equivalent to -1,000,000, and <code>red</code> , <code>yellow</code> , and <code>green</code> are equivalent to integers as described in Tracking Node Health .
text	A text string
timeout	A time duration, specified as a bare number (in which case it is considered to be in seconds) or a number with a unit (<code>ms</code> or <code>msec</code> for milliseconds, <code>us</code> or <code>usec</code> for microseconds, <code>s</code> or <code>sec</code> for seconds, <code>m</code> or <code>min</code> for minutes, <code>h</code> or <code>hr</code> for hours) optionally with whitespace before and/or after the number.
version	Version number (any combination of alphanumeric characters, dots, and dashes, starting with a number).

Scores

Scores are integral to how Pacemaker works. Practically everything from moving a resource to deciding which resource to stop in a degraded cluster is achieved by manipulating scores in some way.

Scores are calculated per resource and node. Any node with a negative score for a resource can't run that resource. The cluster places a resource on the node with the highest score for it.

Score addition and subtraction follow these rules:

- Any value (including `INFINITY`) - `INFINITY` = `-INFINITY`
- `INFINITY` + any value other than `-INFINITY` = `INFINITY`

Note: What if you want to use a score higher than 1,000,000? Typically this possibility arises when someone wants to base the score on some external metric that might go above 1,000,000.

The short answer is you can't.

The long answer is it is sometimes possible work around this limitation creatively. You may be able to set the score to some computed value based on the external metric rather than use the metric directly. For nodes, you can store the metric as a node attribute, and query the attribute when computing the score (possibly as part of a custom resource agent).

2.2.2 Local Options

Most Pacemaker configuration is in the cluster-wide CIB, but some host-local configuration options either are needed at startup (before the CIB is read) or provide per-host overrides of cluster-wide options.

These options are configured as environment variables set when Pacemaker is started, in the format `<NAME>=<VALUE>`. These are typically set in a file whose location varies by OS (most commonly `/etc/sysconfig/pacemaker` or `/etc/default/pacemaker`; this documentation was generated on a system using `/etc/sysconfig/pacemaker`).

Table 2: Local Options

Name	Type	Default	Description
<code>CIB_pam_service</code>	<i>text</i>	login	PAM service to use for remote CIB client authentication (passed to <code>pam_start</code>).
<code>PCMK_logfacility</code>	<i>enumeration</i>	daemon	Enable logging via the system log or journal, using the specified log facility. Messages sent here are of value to all Pacemaker administrators. This can be disabled using <code>none</code> , but that is not recommended. Allowed values: <ul style="list-style-type: none"> • none • daemon • user • local0 • local1 • local2 • local3 • local4 • local5 • local6 • local7
<code>PCMK_logpriority</code>	<i>enumeration</i>	notice	Unless system logging is disabled using <code>PCMK_logfacility=none</code> , messages of the specified log severity and higher will be sent to the system log. The default is appropriate for most installations. Allowed values: <ul style="list-style-type: none"> • emerg • alert • crit • error • warning • notice • info • debug

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
PCMK_logfile	<i>text</i>	<code>/var/log/pacemaker/pacemaker.log</code>	When <code>pacemaker.log</code> is not <code>none</code> , more detailed log messages will be sent to the specified file (in addition to the system log, if enabled). These messages may have extended information, and will include messages of info severity. This log is of more use to developers and advanced system administrators, and when reporting problems. Note: The default is <code>/var/log/pcmk-init.log</code> (inside the container) for bundled container nodes; this would typically be mapped to a different path on the host running the container.
PCMK_logfile_mode	<i>text</i>	0660	Pacemaker will set the permissions on the detail log to this value (see <code>chmod(1)</code>).
PCMK_debug	<i>enumeration</i>	no	Whether to send debug severity messages to the detail log. This may be set for all subsystems (<code>yes</code> or <code>no</code>) or for specific (comma-separated) subsystems. Allowed subsystems are: <ul style="list-style-type: none"> • <code>pacemakerd</code> • <code>pacemaker-attribd</code> • <code>pacemaker-based</code> • <code>pacemaker-controld</code> • <code>pacemaker-execd</code> • <code>pacemaker-fenced</code> • <code>pacemaker-schedulerd</code> Example: <code>PCMK_debug="pacemakerd, pacemaker-execd"</code>
PCMK_stderr	<i>boolean</i>	no	<i>Advanced Use Only:</i> Whether to send daemon log messages to stderr. This would be useful only during troubleshooting, when starting Pacemaker manually on the command line. Setting this option in the configuration file is pointless, since the file is not read when starting Pacemaker manually. However, it can be set directly as an environment variable on the command line.
PCMK_trace_functions	<i>text</i>		<i>Advanced Use Only:</i> Send debug and trace severity messages from these (comma-separated) source code functions to the detail log. Example: <code>PCMK_trace_functions="func1, func2"</code>
PCMK_trace_files	<i>text</i>		<i>Advanced Use Only:</i> Send debug and trace severity messages from all functions in these (comma-separated) source file names to the detail log. Example: <code>PCMK_trace_files="file1.c, file2.c"</code>

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
PCMK_trace_formats	<i>text</i>		<i>Advanced Use Only:</i> Send trace severity messages that are generated by these (comma-separated) format strings in the source code to the detail log. Example: <code>PCMK_trace_formats="Error: %s (%d)"</code>
PCMK_trace_tags	<i>text</i>		<i>Advanced Use Only:</i> Send debug and trace severity messages related to these (comma-separated) resource IDs to the detail log. Example: <code>PCMK_trace_tags="client-ip, dbfs"</code>
PCMK_blackbox	<i>enumeration</i>	no	<i>Advanced Use Only:</i> Enable blackbox logging globally (yes or no) or by subsystem. A blackbox contains a rolling buffer of all logs (of all severities). Blackboxes are stored under <code>/var/lib/pacemaker/blackbox</code> by default, by default, and their contents can be viewed using the <code>qb-blackbox(8)</code> command. The blackbox recorder can be enabled at start using this variable, or at runtime by sending a Pacemaker subsystem daemon process a <code>SIGUSR1</code> or <code>SIGTRAP</code> signal, and disabled by sending <code>SIGUSR2</code> (see <code>kill(1)</code>). The blackbox will be written after a crash, assertion failure, or <code>SIGTRAP</code> signal. See <i>PCMK_debug</i> for allowed subsystems. Example: <code>PCMK_blackbox="pacemakerd, pacemaker-execd"</code>
PCMK_trace_blackbox	<i>enumeration</i>		<i>Advanced Use Only:</i> Write a blackbox whenever the message at the specified function and line is logged. Multiple entries may be comma-separated. Example: <code>PCMK_trace_blackbox="remote.c:144,remote.c:149"</code>
PCMK_node_start_state	<i>enumeration</i>	default	By default, the local host will join the cluster in an online or standby state when Pacemaker first starts depending on whether it was previously put into standby mode. If this variable is set to standby or online , it will force the local host to join in the specified state.
PCMK_node_action_limit	<i>nonnegative integer</i>		If set, this overrides the <i>node-action-limit</i> cluster option on this node to specify the maximum number of jobs that can be scheduled on this node (or 0 to use twice the number of CPU cores).
PCMK_fail_fast	<i>boolean</i>	no	By default, if a Pacemaker subsystem crashes, the main <code>pacemakerd</code> process will attempt to restart it. If this variable is set to yes , <code>pacemakerd</code> will panic the local host instead.

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
PCMK_panic_action	<i>enumeration</i>	reboot	Pacemaker will panic the local host under certain conditions. By default, this means rebooting the host. This variable can change that behavior: if <code>crash</code> , trigger a kernel crash (useful if you want a kernel dump to investigate); if <code>sync-reboot</code> or <code>sync-crash</code> , synchronize filesystems before rebooting the host or triggering a kernel crash. The sync values are more likely to preserve log messages, but with the risk that the host may be left active if the synchronization hangs.
PCMK_remote_address	<i>text</i>		By default, if the <i>Pacemaker Remote</i> service is run on the local node, it will listen for connections on all IP addresses. This may be set to one address to listen on instead, as a resolvable hostname or as a numeric IPv4 or IPv6 address. When resolving names or listening on all addresses, IPv6 will be preferred if available. When listening on an IPv6 address, IPv4 clients will be supported via IPv4-mapped IPv6 addresses. Example: <code>PCMK_remote_address="192.0.2.1"</code>
PCMK_remote_port	<i>port</i>	3121	Use this TCP port number for <i>Pacemaker Remote</i> node connections. This value must be the same on all nodes.
PCMK_ca_file	<i>text</i>		The location of a file containing trusted Certificate Authorities, used to verify client or server certificates. This file must be in PEM format and must be readable by Pacemaker daemons (that is, it must allow read permissions to either the <code>hacluster</code> user or the <code>haclient</code> group). If set, along with <i>PCMK_key_file</i> and <i>PCMK_cert_file</i> , X509 authentication will be enabled for <i>Pacemaker Remote</i> and remote CIB connections. Example: <code>PCMK_ca_file="/etc/pacemaker/ca.cert.pem"</code>
PCMK_cert_file	<i>text</i>		The location of a file containing the signed certificate for the server side of the connection. This file must be in PEM format and must be readable by Pacemaker daemons (that is, it must allow read permissions to either the <code>hacluster</code> user or the <code>haclient</code> group). If set, along with <i>PCMK_ca_file</i> and <i>PCMK_key_file</i> , X509 authentication will be enabled for <i>Pacemaker Remote</i> and remote CIB connections. Example: <code>PCMK_cert_file="/etc/pacemaker/server.cert.pem"</code>

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
PCMK_crl_file	<i>text</i>		The location of a Certificate Revocation List file, in PEM format. This setting is optional for X509 authentication. Example: <code>PCMK_crl_file="/etc/pacemaker/crl.pem"</code>
PCMK_key_file	<i>text</i>		The location of a file containing the private key for the matching <i>PCMK_cert_file</i> , in PEM format. This file must be readable by Pacemaker daemons (that is, it must allow read permissions to either the <code>hacluster</code> user or the <code>haclient</code> group). If set, along with <i>PCMK_ca_file</i> and <i>PCMK_cert_file</i> , X509 authentication will be enabled for <i>Pacemaker Remote</i> and remote CIB connections. Example: <code>PCMK_key_file="/etc/pacemaker/server.key.pem"</code>
PCMK_authkey_location	<i>text</i>	<code>/etc/pacemaker/authkey</code>	As an alternative to using X509 authentication for <i>Pacemaker Remote</i> connections, use the contents of this file as the authorization key. This file must be readable by Pacemaker daemons (that is, it must allow read permissions to either the <code>hacluster</code> user or the <code>haclient</code> group), and its contents must be identical on all nodes. This is an alternative to using X509 certificates.

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
PCMK_remote_pid1	<i>enumeration</i>	default	<p><i>Advanced Use Only:</i> When a bundle resource's <code>run-command</code> option is left to default, <i>Pacemaker Remote</i> runs as PID 1 in the bundle's containers. When it does so, it loads environment variables from the container's <code>/etc/pacemaker/pcmk-init.env</code> and performs the PID 1 responsibility of reaping dead subprocesses.</p> <p>This option controls whether those actions are performed when Pacemaker Remote is not running as PID 1. It is intended primarily for developer testing but can be useful when <code>run-command</code> is set to a separate, custom PID 1 process that launches Pacemaker Remote.</p> <ul style="list-style-type: none"> full: Pacemaker Remote loads environment variables from <code>/etc/pacemaker/pcmk-init.env</code> and reaps dead subprocesses. vars: Pacemaker Remote loads environment variables from <code>/etc/pacemaker/pcmk-init.env</code> but does not reap dead subprocesses. default: Pacemaker Remote performs neither action. <p>If Pacemaker Remote is running as PID 1, this option is ignored, and the behavior is the same as for <code>full</code>.</p>
PCMK_tls_priorities	<i>text</i>	NORMAL	<p><i>Advanced Use Only:</i> These GnuTLS cipher priorities will be used for TLS connections (whether for <i>Pacemaker Remote</i> connections or remote CIB access, when enabled). Pacemaker will append <code>:+ANON-DH</code> for remote CIB access and <code>:+DHE-PSK:+PSK</code> for Pacemaker Remote connections, as they are required for the respective functionality.</p> <p>Example: <code>PCMK_tls_priorities="SECURE128:+SECURE192"</code></p>
PCMK_dh_max_bits	<i>nonnegative integer</i>	0 (no maximum)	<p><i>Advanced Use Only:</i> Set an upper bound on the bit length of the prime number generated for Diffie-Hellman parameters needed by TLS connections. The default is no maximum.</p> <p>The server (<i>Pacemaker Remote</i> daemon, or CIB manager configured to accept remote clients) will use this value to provide a ceiling for the value recommended by the GnuTLS library. The library will only accept a limited number of specific values, which vary by library version, so setting these is recommended only when required for compatibility with specific client versions.</p> <p>Clients do not use <code>PCMK_dh_max_bits</code>.</p>

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
PCMK_ipc_type	<i>enumeration</i>	shared-mem	<i>Advanced Use Only:</i> Force use of a particular IPC method. Allowed values: <ul style="list-style-type: none"> • <code>shared-mem</code> • <code>socket</code> • <code>posix</code> • <code>sysv</code>
PCMK_ipc_buffer	<i>nonnegative integer</i>	131072	<i>Advanced Use Only:</i> Specify an IPC buffer size in bytes. This can be useful when connecting to large clusters that result in messages exceeding the default size (which will also result in log messages referencing this variable).
PCMK_cluster_type	<i>enumeration</i>	corosync	<i>Advanced Use Only:</i> Specify the cluster layer to be used. If unset, Pacemaker will detect and use a supported cluster layer, if available. Currently, " <code>corosync</code> " is the only supported cluster layer. If multiple layers are supported in the future, this will allow overriding Pacemaker's automatic detection to select a specific one.
PCMK_schema_directory	<i>text</i>	<code>/usr/share/pacemaker</code>	<i>Advanced Use Only:</i> Specify an alternate location for RNG schemas and XSL transforms.
PCMK_remote_schema_directory	<i>text</i>	<code>/var/lib/pacemaker</code>	<i>Advanced Use Only:</i> Specify an alternate location on <i>Pacemaker Remote</i> nodes for storing newer RNG schemas and XSL transforms fetched from the cluster.
PCMK_valgrind_enabled	<i>enumeration</i>	no	<i>Advanced Use Only:</i> Whether subsystem daemons should be run under <code>valgrind</code> . Allowed values are the same as for <code>PCMK_debug</code> .
PCMK_callgrind_enabled	<i>enumeration</i>	no	<i>Advanced Use Only:</i> Whether subsystem daemons should be run under <code>valgrind</code> with the <code>callgrind</code> tool enabled. Allowed values are the same as for <code>PCMK_debug</code> .
SBD_SYNC_RESOURCE_STARTUP	<i>boolean</i>		If true, <code>pacemakerd</code> waits for a ping from <code>sbd</code> during startup before starting other Pacemaker daemons, and during shutdown after stopping other Pacemaker daemons but before exiting. Default value is set based on the <code>--with-sbd-sync-default</code> configure script option.
SBD_WATCHDOG_TIMEOUT	<i>duration</i>		If the <code>stonith-watchdog-timeout</code> cluster property is set to a negative or invalid value, use double this value as the default if positive, or use 0 as the default otherwise. This value must be greater than the value of <code>stonith-watchdog-timeout</code> if both are set.

Continued on next page

Table 2 – continued from previous page

Name	Type	Default	Description
VAL-GRIND_OPTS	<i>text</i>		<i>Advanced Use Only:</i> Pass these options to valgrind, when enabled (see <code>valgrind(1)</code>). " <code>--vgdb=no</code> " should usually be specified because <code>pacemaker-execd</code> can lower privileges when executing commands, which would otherwise leave a bunch of unremovable files in <code>/tmp</code> .

2.3 Cluster-Wide Configuration

2.3.1 Configuration Layout

The cluster is defined by the Cluster Information Base (CIB), which uses XML notation. The simplest CIB, an empty one, looks like this:

An empty configuration

```
<cib crm_feature_set="3.6.0" validate-with="pacemaker-3.5" epoch="1" num_updates="0" admin_epoch="0">
  <configuration>
    <crm_config/>
    <nodes/>
    <resources/>
    <constraints/>
  </configuration>
  <status/>
</cib>
```

The empty configuration above contains the major sections that make up a CIB:

- **cib:** The entire CIB is enclosed with a `cib` element. Certain fundamental settings are defined as attributes of this element.
 - **configuration:** This section – the primary focus of this document – contains traditional configuration information such as what resources the cluster serves and the relationships among them.
 - * **crm_config:** cluster-wide configuration options
 - * **nodes:** the machines that host the cluster
 - * **resources:** the services run by the cluster
 - * **constraints:** indications of how resources should be placed
 - **status:** This section contains the history of each resource on each node. Based on this data, the cluster can construct the complete current state of the cluster. The authoritative source for this section is the local executor (`pacemaker-execd` process) on each cluster node, and the cluster will occasionally repopulate the entire section. For this reason, it is never written to disk, and administrators are advised against modifying it in any way.

In this document, configuration settings will be described as properties or options based on how they are defined in the CIB:

- Properties are XML attributes of an XML element.

- Options are name-value pairs expressed as `nvpair` child elements of an XML element.

Normally, you will use command-line tools that abstract the XML, so the distinction will be unimportant; both properties and options are cluster settings you can tweak.

Options can appear within four types of enclosing elements:

- `cluster_property_set`
- `instance_attributes`
- `meta_attributes`
- `utilization`

We will refer to a set of options and its enclosing element as a *block*.

Table 3: Properties of an Option Block's Enclosing Element

Name	Type	Default	Description
<code>id</code>	<i>id</i>		A unique name for the block (required)
<code>score</code>	<i>score</i>	0	Priority with which to process the block

Each block may optionally contain a *rule*.

2.3.2 Option Precedence

This subsection describes the precedence of options within a set of blocks and within a single block.

Options are processed as follows:

- All option blocks of a given type are processed in order of their `score` attribute, from highest to lowest. For `cluster_property_set`, if there is a block whose enclosing element has `id="cib-bootstrap-options"`, then that block is always processed first regardless of score.
- If a block contains a rule that evaluates to false, that block is skipped.
- Within a block, options are processed in order from first to last.
- The first value found for a given option is applied, and the rest are ignored.

Note that this means it is pointless to configure the same option twice in a single block, because occurrences after the first one would be ignored.

For example, in the following configuration snippet, the `no-quorum-policy` value `demote` is applied. `property-set2` has a higher score than `property-set1`, so it's processed first. There are no rules in this snippet, so both sets are processed. Within `property-set2`, the value `demote` appears first, so the later value `freeze` is ignored. We've already found a value for `no-quorum-policy` before we begin processing `property-set1`, so its value `stop` is ignored.

```
<cluster_property_set id="property-set1" score="500">
  <nvpair id="no-quorum-policy1" name="no-quorum-policy" value="stop"/>
</cluster_property_set>
<cluster_property_set id="property-set2" score="1000">
  <nvpair id="no-quorum-policy2a" name="no-quorum-policy" value="demote"/>
  <nvpair id="no-quorum-policy2b" name="no-quorum-policy" value="freeze"/>
</cluster_property_set>
```

2.3.3 CIB Properties

Certain settings are defined by CIB properties (that is, attributes of the `cib` tag) rather than with the rest of the cluster configuration in the `configuration` section.

The reason is simply a matter of parsing. These options are used by the configuration database which is, by design, mostly ignorant of the content it holds. So the decision was made to place them in an easy-to-find location.

Table 4: CIB Properties

Name	Type	Default	Description
<code>admin_epoch</code>	<i>nonnegative integer</i>	0	When a node joins the cluster, the cluster asks the node with the highest (<code>admin_epoch</code> , <code>epoch</code> , <code>num_updates</code>) tuple to replace the configuration on all the nodes – which makes setting them correctly very important. <code>admin_epoch</code> is never modified by the cluster; you can use this to make the configurations on any inactive nodes obsolete.
<code>epoch</code>	<i>nonnegative integer</i>	0	The cluster increments this every time the CIB’s configuration section is updated.
<code>num_updates</code>	<i>nonnegative integer</i>	0	The cluster increments this every time the CIB’s configuration or status sections are updated, and resets it to 0 when epoch changes.
<code>validate-with</code>	<i>enumeration</i>		Determines the type of XML validation that will be done on the configuration. Allowed values are <code>none</code> (in which case the cluster will not require that updates conform to expected syntax) and the base names of schema files installed on the local machine (for example, “ <code>pacemaker-3.9</code> ”)
<code>remote-tls-port</code>	<i>port</i>		If set, the CIB manager will listen for anonymously encrypted remote connections on this port, to allow CIB administration from hosts not in the cluster. No key is used, so this should be used only on a protected network where man-in-the-middle attacks can be avoided.
<code>remote-clear-port</code>	<i>port</i>		If set to a TCP port number, the CIB manager will listen for remote connections on this port, to allow for CIB administration from hosts not in the cluster. No encryption is used, so this should be used only on a protected network.
<code>cib-last-written</code>	<i>date/time</i>		Indicates when the configuration was last written to disk. Maintained by the cluster; for informational purposes only.
<code>have-quorum</code>	<i>boolean</i>		Indicates whether the cluster has quorum. If false, the cluster’s response is determined by <code>no-quorum-policy</code> (see below). Maintained by the cluster.
<code>dc-uuid</code>	<i>text</i>		Node ID of the cluster’s current designated controller (DC). Used and maintained by the cluster.

Continued on next page

Table 4 – continued from previous page

Name	Type	Default	Description
execution-date	<i>epoch time</i>		Time to use when evaluating rules.

2.3.4 Cluster Options

Cluster options, as you might expect, control how the cluster behaves when confronted with various situations.

They are grouped into sets within the `crm_config` section. In advanced configurations, there may be more than one set. (This will be described later in the chapter on *Rules* where we will show how to have the cluster use different sets of options during working hours than during weekends.) For now, we will describe the simple case where each option is present at most once.

You can obtain an up-to-date list of cluster options, including their default values, by running the `man pacemaker-schedulerd` and `man pacemaker-controld` commands.

Table 5: Cluster Options

Name	Type	Default	Description
cluster-name	<i>text</i>		An (optional) name for the cluster as a whole. This is mostly for users' convenience for use as desired in administration, but can be used in the Pacemaker configuration in <i>Rules</i> (as the <code>#cluster-name</code> <i>node attribute</i>). It may also be used by higher-level tools when displaying cluster information, and by certain resource agents (for example, the <code>ocf:heartbeat:GFS2</code> agent stores the cluster name in filesystem meta-data).
dc-version	<i>version</i>	<i>detected</i>	Version of Pacemaker on the cluster's designated controller (DC). Maintained by the cluster, and intended for diagnostic purposes.
cluster-infrastructure	<i>text</i>	<i>detected</i>	The messaging layer with which Pacemaker is currently running. Maintained by the cluster, and intended for informational and diagnostic purposes.

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
no-quorum-policy	<i>enumeration</i>	stop	What to do when the cluster does not have quorum. Allowed values: <ul style="list-style-type: none"> • ignore: continue all resource management • freeze: continue resource management, but don't recover resources from nodes not in the affected partition • stop: stop all resources in the affected cluster partition • demote: demote promotable resources and stop all other resources in the affected cluster partition (<i>since 2.0.5</i>) • fence: fence all nodes in the affected cluster partition (<i>since 2.1.9</i>) • suicide: same as fence (<i>deprecated since 2.1.9</i>)
batch-limit	<i>integer</i>	0	The maximum number of actions that the cluster may execute in parallel across all nodes. The ideal value will depend on the speed and load of your network and cluster nodes. If zero, the cluster will impose a dynamically calculated limit only when any node has high load. If -1, the cluster will not impose any limit.
migration-limit	<i>integer</i>	-1	The number of <i>live migration</i> actions that the cluster is allowed to execute in parallel on a node. A value of -1 means unlimited.
load-threshold	<i>percentage</i>	80%	Maximum amount of system load that should be used by cluster nodes. The cluster will slow down its recovery process when the amount of system resources used (currently CPU) approaches this limit.
node-action-limit	<i>integer</i>	0	Maximum number of jobs that can be scheduled per node. If nonpositive or invalid, double the number of cores is used as the maximum number of jobs per node. <i>PCMK_node_action_limit</i> overrides this option on a per-node basis.
symmetric-cluster	<i>boolean</i>	true	If true, resources can run on any node by default. If false, a resource is allowed to run on a node only if a <i>location constraint</i> enables it.
stop-all-resources	<i>boolean</i>	false	Whether all resources should be disallowed from running (can be useful during maintenance or troubleshooting)
stop-orphan-resources	<i>boolean</i>	true	Whether resources that have been deleted from the configuration should be stopped. This value takes precedence over <i>is-managed</i> (that is, even unmanaged resources will be stopped when orphaned if this value is true).

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
stop-orphan-actions	<i>boolean</i>	true	Whether recurring <i>operations</i> that have been deleted from the configuration should be cancelled
start-failure-is-fatal	<i>boolean</i>	true	Whether a failure to start a resource on a particular node prevents further start attempts on that node. If false , the cluster will decide whether the node is still eligible based on the resource's current failure count and <i>migration-threshold</i> .
enable-startup-probes	<i>boolean</i>	true	Whether the cluster should check the pre-existing state of resources when the cluster starts
maintenance-mode	<i>boolean</i>	false	If true, the cluster will not start or stop any resource in the cluster, and any recurring operations (except those specifying role as Stopped) will be paused. If true, this overrides the <i>maintenance</i> node attribute, <i>is-managed</i> and <i>maintenance</i> resource meta-attributes, and <i>enabled</i> operation meta-attribute.
stonith-enabled	<i>boolean</i>	true	Whether the cluster is allowed to fence nodes (for example, failed nodes and nodes with resources that can't be stopped). If true, at least one fence device must be configured before resources are allowed to run. If false, unresponsive nodes are immediately assumed to be running no resources, and resource recovery on online nodes starts without any further protection (which can mean <i>data loss</i> if the unresponsive node still accesses shared storage, for example). See also the <i>requires</i> resource meta-attribute. This option applies only to fencing scheduled by the cluster, not to requests initiated externally (such as with the <code>stonith_admin</code> command-line tool).
stonith-action	<i>enumeration</i>	reboot	Action the cluster should send to the fence agent when a node must be fenced. Allowed values are reboot and off .
stonith-timeout	<i>duration</i>	60s	How long to wait for on , off , and reboot fence actions to complete by default.
stonith-max-attempts	<i>score</i>	10	How many times fencing can fail for a target before the cluster will no longer immediately re-attempt it. Any value below 1 will be ignored, and the default will be used instead.

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
have-watchdog	<i>boolean</i>	<i>detected</i>	Whether watchdog integration is enabled. This is set automatically by the cluster according to whether SBD is detected to be in use. User-configured values are ignored. The value <i>true</i> is meaningful if diskless SBD is used and <i>stonith-watchdog-timeout</i> is nonzero. In that case, if fencing is required, watchdog-based self-fencing will be performed via SBD without requiring a fencing resource explicitly configured.
stonith-watchdog-timeout	<i>timeout</i>	0	<p>If nonzero, and the cluster detects have-watchdog as true, then watchdog-based self-fencing will be performed via SBD when fencing is required.</p> <p>If this is set to a positive value, lost nodes are assumed to achieve self-fencing within this much time.</p> <p>This does not require a fencing resource to be explicitly configured, though a <code>fence_watchdog</code> resource can be configured, to limit use to specific nodes.</p> <p>If this is set to 0 (the default), the cluster will never assume watchdog-based self-fencing.</p> <p>If this is set to a negative value, the cluster will use twice the local value of the <code>SBD_WATCHDOG_TIMEOUT</code> environment variable if that is positive, or otherwise treat this as 0.</p> <p>Warning: When used, this timeout must be larger than <code>SBD_WATCHDOG_TIMEOUT</code> on all nodes that use watchdog-based SBD, and Pacemaker will refuse to start on any of those nodes where this is not true for the local value or SBD is not active. When this is set to a negative value, <code>SBD_WATCHDOG_TIMEOUT</code> must be set to the same value on all nodes that use SBD, otherwise data corruption or loss could occur.</p>
concurrent-fencing	<i>boolean</i>	false	Whether the cluster is allowed to initiate multiple fence actions concurrently. Fence actions initiated externally, such as via the <code>stonith_admin</code> tool or an application such as DLM, or by the fencer itself such as recurring device monitors and <code>status</code> and <code>list</code> commands, are not limited by this option.

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
fence-reaction	<i>enumeration</i>	stop	How should a cluster node react if notified of its own fencing? A cluster node may receive notification of a “succeeded” fencing that targeted it if fencing is misconfigured, or if fabric fencing is in use that doesn’t cut cluster communication. Allowed values are stop to attempt to immediately stop Pacemaker and stay stopped, or panic to attempt to immediately reboot the local node, falling back to stop on failure. The default is likely to be changed to panic in a future release. (<i>since 2.0.3</i>)
priority-fencing-delay	<i>duration</i>	0	Apply this delay to any fencing targeting the lost nodes with the highest total resource priority in case we don’t have the majority of the nodes in our cluster partition, so that the more significant nodes potentially win any fencing match (especially meaningful in a split-brain of a 2-node cluster). A promoted resource instance takes the resource’s priority plus 1 if the resource’s priority is not 0. Any static or random delays introduced by <code>pcm_delay_base</code> and <code>pcm_delay_max</code> configured for the corresponding fencing resources will be added to this delay. This delay should be significantly greater than (safely twice) the maximum delay from those parameters. (<i>since 2.0.4</i>)
node-pending-timeout	<i>duration</i>	0	Fence nodes that do not join the controller process group within this much time after joining the cluster, to allow the cluster to continue managing resources. A value of 0 means never fence pending nodes. Setting the value to 2h means fence nodes after 2 hours. (<i>since 2.1.7</i>)
cluster-delay	<i>duration</i>	60s	If the DC requires an action to be executed on another node, it will consider the action failed if it does not get a response from the other node within this time (beyond the action’s own timeout). The ideal value will depend on the speed and load of your network and cluster nodes.
dc-deadtime	<i>duration</i>	20s	How long to wait for a response from other nodes when electing a DC. The ideal value will depend on the speed and load of your network and cluster nodes.
cluster-ipc-limit	<i>nonnegative integer</i>	500	The maximum IPC message backlog before one cluster daemon will disconnect another. This is of use in large clusters, for which a good value is the number of resources in the cluster multiplied by the number of nodes. The default of 500 is also the minimum. Raise this if you see “Evicting client” log messages for cluster daemon process IDs.

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
pe-error-series-max	<i>integer</i>	-1	The number of scheduler inputs resulting in errors to save. These inputs can be helpful during troubleshooting and when reporting issues. A negative value means save all inputs, and 0 means save none.
pe-warn-series-max	<i>integer</i>	5000	The number of scheduler inputs resulting in warnings to save. These inputs can be helpful during troubleshooting and when reporting issues. A negative value means save all inputs, and 0 means save none.
pe-input-series-max	<i>integer</i>	4000	The number of “normal” scheduler inputs to save. These inputs can be helpful during troubleshooting and when reporting issues. A negative value means save all inputs, and 0 means save none.
enable-acl	<i>boolean</i>	false	Whether <i>access control lists</i> should be used to authorize CIB modifications
placement-strategy	<i>enumeration</i>	default	How the cluster should assign resources to nodes (see <i>Utilization and Placement Strategy</i>). Allowed values are default , utilization , balanced , and minimal .
node-health-strategy	<i>enumeration</i>	none	How the cluster should react to <i>node health</i> attributes. Allowed values are none , migrate-on-red , only-green , progressive , and custom .
node-health-base	<i>score</i>	0	The base health score assigned to a node. Only used when node-health-strategy is progressive .
node-health-green	<i>score</i>	0	The score to use for a node health attribute whose value is green . Only used when node-health-strategy is progressive or custom .
node-health-yellow	<i>score</i>	0	The score to use for a node health attribute whose value is yellow . Only used when node-health-strategy is progressive or custom .
node-health-red	<i>score</i>	-INFINITY	The score to use for a node health attribute whose value is red . Only used when node-health-strategy is progressive or custom .

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
cluster-recheck-interval	<i>duration</i>	15min	<p>Pacemaker is primarily event-driven, and looks ahead to know when to recheck the cluster for failure-timeout settings and most time-based rules (<i>since 2.0.3</i>). However, it will also recheck the cluster after this amount of inactivity. This has three main effects:</p> <ul style="list-style-type: none"> • <i>Rules</i> using <code>date_spec</code> are guaranteed to be checked only this often. • If <i>fencing</i> fails enough to reach <i>stonith-max-attempts</i>, attempts will begin again after at most this time. • It serves as a fail-safe in case of certain scheduler bugs. If the scheduler incorrectly determines only some of the actions needed to react to a particular event, it will often correctly determine the rest after at most this time. <p>A value of 0 disables this polling.</p>
shutdown-lock	<i>boolean</i>	false	<p>The default of false allows active resources to be recovered elsewhere when their node is cleanly shut down, which is what the vast majority of users will want. However, some users prefer to make resources highly available only for failures, with no recovery for clean shutdowns. If this option is true, resources active on a node when it is cleanly shut down are kept “locked” to that node (not allowed to run elsewhere) until they start again on that node after it rejoins (or for at most <code>shutdown-lock-limit</code>, if set). Stonith resources and Pacemaker Remote connections are never locked. Clone and bundle instances and the promoted role of promotable clones are currently never locked, though support could be added in a future release. Locks may be manually cleared using the <code>--refresh</code> option of <code>crm_resource</code> (both the resource and node must be specified; this works with remote nodes if their connection resource’s <code>target-role</code> is set to <code>Stopped</code>, but not if Pacemaker Remote is stopped on the remote node without disabling the connection resource). (<i>since 2.0.4</i>)</p>
shutdown-lock-limit	<i>duration</i>	0	<p>If <code>shutdown-lock</code> is true, and this is set to a nonzero time duration, locked resources will be allowed to start after this much time has passed since the node shutdown was initiated, even if the node has not rejoined. (This works with remote nodes only if their connection resource’s <code>target-role</code> is set to <code>Stopped</code>.) (<i>since 2.0.4</i>)</p>

Continued on next page

Table 5 – continued from previous page

Name	Type	Default	Description
startup-fencing	<i>boolean</i>	true	<i>Advanced Use Only:</i> Whether the cluster should fence unseen nodes at start-up. Setting this to false is unsafe, because the unseen nodes could be active and running resources but unreachable. <code>dc-deadtime</code> acts as a grace period before this fencing, since a DC must be elected to schedule fencing.
election-timeout	<i>duration</i>	2min	<i>Advanced Use Only:</i> If a winner is not declared within this much time of starting an election, the node that initiated the election will declare itself the winner.
shutdown-escalation	<i>duration</i>	20min	<i>Advanced Use Only:</i> The controller will exit immediately if a shutdown does not complete within this much time.
join-integration-timeout	<i>duration</i>	3min	<i>Advanced Use Only:</i> If you need to adjust this value, it probably indicates the presence of a bug.
join-finalization-timeout	<i>duration</i>	30min	<i>Advanced Use Only:</i> If you need to adjust this value, it probably indicates the presence of a bug.
transition-delay	<i>duration</i>	0s	<i>Advanced Use Only:</i> Delay cluster recovery for the configured interval to allow for additional or related events to occur. This can be useful if your configuration is sensitive to the order in which ping updates arrive. Enabling this option will slow down cluster recovery under all conditions.

2.4 Nodes

Pacemaker supports two basic types of nodes: *cluster nodes* and *Pacemaker Remote nodes*.

2.4.1 Cluster nodes

Cluster nodes run Corosync and all Pacemaker components. They may run cluster resources, run all Pacemaker command-line tools, execute fencing actions, count toward cluster quorum, and serve as the cluster's Designated Controller (DC).

Every cluster must have at least one cluster node. Scalability is limited by the cluster layer to around 32 cluster nodes.

Host Clock Considerations

In general, Pacemaker does not rely on time or time zones being synchronized across nodes. However, if the configuration uses date/time-based *rules*, synchronization is a good idea, otherwise the rules will evaluate differently depending on which node is the Designated Controller (DC). Also, synchronization is greatly helpful when comparing logs across multiple nodes for problem investigation.

If a node's clock jumps forward, you may see relatively minor issues such as various timeouts suddenly being considered expired.

If a node's clock jumps backward, more serious problems may occur, so this should be avoided. If the host clock is adjusted at boot, and Pacemaker is enabled at boot, Pacemaker's start should be ordered after the clock adjustment. When run under `systemd`, Pacemaker will automatically order itself after `time-sync.target`. However, depending on the local setup, you may need to enable an additional service (for example, `chronyd-wait.service`) for that to be effective, or write your own workaround (for example, see the discussion on [systemd issue#5097](#)).

2.4.2 Pacemaker Remote nodes

Pacemaker Remote nodes do not run Corosync or the usual Pacemaker components. Instead, they run only the *remote executor* (`pacemaker-remoted`), which waits for Pacemaker on a cluster node to give it instructions.

They may run cluster resources and most command-line tools, but cannot perform other functions of full cluster nodes such as fencing execution, quorum voting, or DC eligibility.

There is no hard limit on the number of Pacemaker Remote nodes.

Note: *Remote* in this document has nothing to do with physical proximity and instead refers to the node not being a member of the underlying Corosync cluster. Pacemaker Remote nodes are subject to the same latency requirements as cluster nodes, which means they are typically in the same data center.

There are three types of Pacemaker Remote nodes:

- A *remote node* boots outside Pacemaker control, and is typically a physical host. The connection to the remote node is managed as a *special type of resource* configured by the user.
- A *guest node* is a virtual machine or container configured to run Pacemaker's remote executor when launched, and is launched and managed by the cluster as a standard resource configured by the user with *special options*.
- A *bundle node* is a guest node created for a container that is launched and managed by the cluster as part of a *bundle* resource configured by the user.

Note: It is important to distinguish the various roles a virtual machine can serve in Pacemaker clusters:

- A virtual machine can run the full cluster stack, in which case it is a cluster node and is not itself managed by the cluster.
 - A virtual machine can be managed by the cluster as a simple resource, without the cluster having any awareness of the services running within it. The virtual machine is *opaque* to the cluster.
 - A virtual machine can be a guest node, allowing the cluster to manage both the virtual machine and resources running within it. The virtual machine is *transparent* to the cluster.
-

2.4.3 Defining a Node

Each cluster node will have an entry in the `nodes` section containing at least an ID and a name. A cluster node's ID is defined by the cluster layer (Corosync).

Example Corosync cluster node entry

```
<node id="101" uname="pcmk-1"/>
```

Pacemaker Remote nodes are defined by a resource in the **resources** section. Remote nodes and guest nodes may optionally have an entry in the **nodes** section, primarily for permanent *node attributes*.

Normally, the user should let the cluster populate the **nodes** section automatically.

Where Pacemaker Gets the Node Name

The name that Pacemaker uses for a node in the configuration does not have to be the same as its local hostname. Pacemaker uses the following for a cluster node's name, in order of most preferred first:

- The value of **name** in the **nodelist** section of **corosync.conf** (**nodeid** must also be explicitly set there in order for Pacemaker to associate the name with the node)
- The value of **ring0_addr** in the **nodelist** section of **corosync.conf**
- The local hostname (value of **uname -n**)

A Pacemaker Remote node's name is defined in its resource configuration.

If the cluster is running, the **crm_node -n** command will display the local node's name as used by the cluster.

If a Corosync **nodelist** is used, **crm_node --name-for-id** with a Corosync node ID will display the name used by the node with the given Corosync **nodeid**, for example:

```
crm_node --name-for-id 2
```

2.4.4 Quorum-only Nodes

One popular cluster design uses an even number of cluster nodes (often 2), with an additional lightweight host that contributes to providing quorum but cannot run resources.

With Pacemaker, this can be achieved in either of two ways:

- When Corosync is used as the underlying cluster layer, the lightweight host can run **qdevice** instead of Corosync and Pacemaker.
- The lightweight host can be configured as a Pacemaker cluster node, and a *location constraint* can be configured for the node with **score** set to **-INFINITY**, **rsc-pattern** set to **.***, and **resource-discovery** set to **never**.

2.4.5 Node Attributes

Pacemaker allows node-specific values to be specified using *node attributes*. A node attribute has a name, and may have a distinct value for each node.

Node attributes come in two types, *permanent* and *transient*. Permanent node attributes are kept within the **node** entry, and keep their values even if the cluster restarts on a node. Transient node attributes are kept in the CIB's **status** section, and go away when the cluster stops on the node.

While certain node attributes have specific meanings to the cluster, they are mainly intended to allow administrators and resource agents to track any information desired.

For example, an administrator might choose to define node attributes for how much RAM and disk space each node has, which OS each uses, or which server room rack each node is in.

Users can configure *Rules* that use node attributes to affect where resources are placed.

Setting and querying node attributes

Node attributes can be set and queried using the `crm_attribute` and `attrd_updater` commands, so that the user does not have to deal with XML configuration directly.

Here is an example command to set a permanent node attribute, and the XML configuration that would be generated:

Result of using `crm_attribute` to specify which kernel `pcmk-1` is running

```
# crm_attribute --type nodes --node pcmk-1 --name kernel --update $(uname -r)

<node id="1" uname="pcmk-1">
  <instance_attributes id="nodes-1-attributes">
    <nvpair id="nodes-1-kernel" name="kernel" value="3.10.0-862.14.4.el7.x86_64"/>
  </instance_attributes>
</node>
```

To read back the value that was just set:

```
# crm_attribute --type nodes --node pcmk-1 --name kernel --query
scope=nodes name=kernel value=3.10.0-862.14.4.el7.x86_64
```

The `--type nodes` indicates that this is a permanent node attribute; `--type status` would indicate a transient node attribute.

Warning: Attribute values with newline or tab characters are currently displayed with newlines as `"\n"` and tabs as `"\t"`, when `crm_attribute` or `attrd_updater` query commands use `--output-as=text` or leave `--output-as` unspecified:

```
# crm_attribute -N node1 -n test_attr -v "$(echo -e "a\nb\tc")" -t status
# crm_attribute -N node1 -n test_attr --query -t status
scope=status name=test_attr value=a\nb\tc
```

This format is deprecated. In a future release, the values will be displayed with literal whitespace characters:

```
# crm_attribute -N node1 -n test_attr --query -t status
scope=status name=test_attr value=a
b c
```

Users should either avoid attribute values with newlines and tabs, or ensure that they can handle both formats.

However, it's best to use `--output-as=xml` when parsing attribute values from output. Newlines, tabs, and special characters are replaced with XML character references that a conforming XML processor can recognize and convert to literals (*since 2.1.8*):

```
# crm_attribute -N node1 -n test_attr --query -t status --output-as=xml
<pacemaker-result api-version="2.35" request="crm_attribute -N laptop -n test_attr --query -t
↪status --output-as=xml">
  <attribute name="test_attr" value="a&#10;b&#9;c" scope="status"/>
  <status code="0" message="OK"/>
</pacemaker-result>
```

Special node attributes

Certain node attributes have special meaning to the cluster.

Node attribute names beginning with `#` are considered reserved for these special attributes. Some special attributes do not start with `#`, for historical reasons.

Certain special attributes are set automatically by the cluster, should never be modified directly, and can be used only within *Rules*; these are listed under *built-in node attributes*.

For true/false values, the cluster considers a value of “1”, “y”, “yes”, “on”, or “true” (case-insensitively) to be true, “0”, “n”, “no”, “off”, “false”, or unset to be false, and anything else to be an error.

Table 6: Node attributes with special significance

Name	Description
fail-count-*	Attributes whose names start with <code>fail-count-</code> are managed by the cluster to track how many times particular resource operations have failed on this node. These should be queried and cleared via the <code>crm_failcount</code> or <code>crm_resource --cleanup</code> commands rather than directly.
last-failure-*	Attributes whose names start with <code>last-failure-</code> are managed by the cluster to track when particular resource operations have most recently failed on this node. These should be cleared via the <code>crm_failcount</code> or <code>crm_resource --cleanup</code> commands rather than directly.
maintenance	If true, the cluster will not start or stop any resources on this node. Any resources active on the node become unmanaged, and any recurring operations for those resources (except those specifying <code>role as Stopped</code>) will be paused. The <i>maintenance-mode</i> cluster option, if true, overrides this. If this attribute is true, it overrides the <i>is-managed</i> and <i>maintenance</i> meta-attributes of affected resources and <i>enabled</i> meta-attribute for affected recurring actions. Pacemaker should not be restarted on a node that is in single-node maintenance mode.
probe_complete	This is managed by the cluster to detect when nodes need to be probed, and should never be used directly.
resource-discovery-enabled	If the node is a remote node, fencing is enabled, and this attribute is explicitly set to false (unset means true in this case), resource discovery (probes) will not be done on this node. This is highly discouraged; the <code>resource-discovery</code> location constraint property is preferred for this purpose.
shutdown	This is managed by the cluster to orchestrate the shutdown of a node, and should never be used directly.
site-name	If set, this will be used as the value of the <code>#site-name</code> node attribute used in rules. (If not set, the value of the <code>cluster-name</code> cluster option will be used as <code>#site-name</code> instead.)
standby	If true, the node is in standby mode. This is typically set and queried via the <code>crm_standby</code> command rather than directly.
terminate	If the value is true or begins with any nonzero number, the node will be fenced. This is typically set by tools rather than directly.

Continued on next page

Table 6 – continued from previous page

Name	Description
<code>#digests-*</code>	Attributes whose names start with <code>#digests-</code> are managed by the cluster to detect when <i>Unfencing</i> needs to be redone, and should never be used directly.
<code>#node-unfenced</code>	When the node was last unfenced (as seconds since the epoch). This is managed by the cluster and should never be used directly.

2.4.6 Tracking Node Health

A node may be functioning adequately as far as cluster membership is concerned, and yet be “unhealthy” in some respect that makes it an undesirable location for resources. For example, a disk drive may be reporting SMART errors, or the CPU may be highly loaded.

Pacemaker offers a way to automatically move resources off unhealthy nodes.

Node Health Attributes

Pacemaker will treat any node attribute whose name starts with `#health` as an indicator of node health. Node health attributes may have one of the following values:

Table 7: Allowed Values for Node Health Attributes

Value	Intended significance
<code>red</code>	This indicator is unhealthy
<code>yellow</code>	This indicator is close to unhealthy (whether worsening or recovering)
<code>green</code>	This indicator is healthy
<i>integer</i>	A numeric score to apply to all resources on this node (0 or positive is healthy, negative is unhealthy)

Note: A health attribute may technically be transient or permanent, but generally only transient makes sense.

Note: `red`, `yellow`, and `green` function as aliases for particular numeric scores as described later.

Node Health Strategy

Pacemaker assigns a node health score to each node, as the sum of the values of all its node health attributes. This score will be used as a location constraint applied to this node for all resources.

The `node-health-strategy` cluster option controls how Pacemaker responds to changes in node health attributes, and how it translates `red`, `yellow`, and `green` to scores.

Allowed values are:

Table 8: Node Health Strategies

Value	Effect
none	Do not track node health attributes at all.
migrate-on-red	Assign the value of <code>-INFINITY</code> to <code>red</code> , and 0 to <code>yellow</code> and <code>green</code> . This will cause all resources to move off the node if any attribute is <code>red</code> .
only-green	Assign the value of <code>-INFINITY</code> to <code>red</code> and <code>yellow</code> , and 0 to <code>green</code> . This will cause all resources to move off the node if any attribute is <code>red</code> or <code>yellow</code> .
progressive	Assign the value of the <code>node-health-red</code> cluster option to <code>red</code> , the value of <code>node-health-yellow</code> to <code>yellow</code> , and the value of <code>node-health-green</code> to <code>green</code> . Each node is additionally assigned a score of <code>node-health-base</code> (this allows resources to start even if some attributes are <code>yellow</code>). This strategy gives the administrator finer control over how important each value is.
custom	Track node health attributes using the same values as <code>progressive</code> for <code>red</code> , <code>yellow</code> , and <code>green</code> , but do not take them into account. The administrator is expected to implement a policy by defining <i>Rules</i> referencing node health attributes.

Exempting a Resource from Health Restrictions

If you want a resource to be able to run on a node even if its health score would otherwise prevent it, set the resource's `allow-unhealthy-nodes` meta-attribute to `true` (*available since 2.1.3*).

This is particularly useful for node health agents, to allow them to detect when the node becomes healthy again. If you configure a health agent without this setting, then the health agent will be banned from an unhealthy node, and you will have to investigate and clear the health attribute manually once it is healthy to allow resources on the node again.

If you want the meta-attribute to apply to a clone, it must be set on the clone itself, not on the resource being cloned.

Configuring Node Health Agents

Since Pacemaker calculates node health based on node attributes, any method that sets node attributes may be used to measure node health. The most common are resource agents and custom daemons.

Pacemaker provides examples that can be used directly or as a basis for custom code. The `ocf:pacemaker:HealthCPU`, `ocf:pacemaker:HealthIOWait`, and `ocf:pacemaker:HealthSMART` resource agents set node health attributes based on CPU and disk status.

To take advantage of this feature, add the resource to your cluster (generally as a cloned resource with a recurring monitor action, to continually check the health of all nodes). For example:

Example HealthIOWait resource configuration

```

<clone id="resHealthIOWait-clone">
  <primitive class="ocf" id="HealthIOWait" provider="pacemaker" type="HealthIOWait">
    <instance_attributes id="resHealthIOWait-instance_attributes">
      <nvpair id="resHealthIOWait-instance_attributes-red_limit" name="red_limit" value="30"/>
      <nvpair id="resHealthIOWait-instance_attributes-yellow_limit" name="yellow_limit" value="10
↵"/>
    </instance_attributes>
    <operations>
      <op id="resHealthIOWait-monitor-interval-5" interval="5" name="monitor" timeout="5"/>
      <op id="resHealthIOWait-start-interval-0s" interval="0s" name="start" timeout="10s"/>
      <op id="resHealthIOWait-stop-interval-0s" interval="0s" name="stop" timeout="10s"/>
    </operations>
  </primitive>
</clone>

```

The resource agents use `attrd_updater` to set proper status for each node running this resource, as a node attribute whose name starts with `#health` (for `HealthIOWait`, the node attribute is named `#health-iowait`).

When a node is no longer faulty, you can force the cluster to make it available to take resources without waiting for the next monitor, by setting the node health attribute to green. For example:

Force node1 to be marked as healthy

```
# attrd_updater --name "#health-iowait" --update "green" --node "node1"
```

2.5 Resources

A *resource* is a service managed by Pacemaker. The simplest type of resource, a *primitive*, is described in this chapter. More complex forms, such as groups and clones, are described in later chapters.

Every primitive has a *resource agent* that provides Pacemaker a standardized interface for managing the service. This allows Pacemaker to be agnostic about the services it manages. Pacemaker doesn't need to understand how the service works because it relies on the resource agent to do the right thing when asked.

Every resource has a *standard* (also called *class*) specifying the interface that its resource agent follows, and a *type* identifying the specific service being managed.

2.5.1 Resource Standards

Pacemaker can use resource agents complying with these standards, described in more detail below:

- ocf
- lsb
- systemd
- service
- stonith

Support for some standards is controlled by build options and so might not be available in any particular build of Pacemaker. The command `crm_resource --list-standards` will show which standards are supported by the local build.

Open Cluster Framework

The Open Cluster Framework (OCF) Resource Agent API is a ClusterLabs standard for managing services. It is the most preferred since it is specifically designed for use in a Pacemaker cluster.

OCF agents are scripts that support a variety of actions including `start`, `stop`, and `monitor`. They may accept parameters, making them more flexible than other standards. The number and purpose of parameters is left to the agent, which advertises them via the `meta-data` action.

Unlike other standards, OCF agents have a *provider* as well as a standard and type.

For more information, see the “Resource Agents” chapter of *Pacemaker Administration* and the [OCF standard](#).

Systemd

Most Linux distributions use `Systemd` for system initialization and service management. *Unit files* specify how to manage services and are usually provided by the distribution.

Pacemaker can manage `systemd` units of type `service`, `socket`, `mount`, `timer`, or `path`. Simply create a resource with `systemd` as the resource standard and the unit file name as the resource type. Do *not* run `systemctl enable` on the unit.

Important: Make sure that any `systemd` services to be controlled by the cluster are *not* enabled to start at boot.

Linux Standard Base

LSB resource agents, also known as *SysV-style*, are scripts that provide `start`, `stop`, and `status` actions for a service.

They are provided by some operating system distributions. If a full path is not given, they are assumed to be located in a directory specified when your Pacemaker software was built (usually `/etc/init.d`).

In order to be used with Pacemaker, they must conform to the [LSB specification](#) as it relates to `init` scripts.

Warning: Some *LSB* scripts do not fully comply with the standard. For details on how to check whether your script is *LSB*-compatible, see the “Resource Agents” chapter of *Pacemaker Administration*. Common problems include:

- Not implementing the `status` action
- Not observing the correct exit status codes
- Starting a started resource returns an error
- Stopping a stopped resource returns an error

Important: Make sure the host is *not* configured to start any *LSB* services at boot that will be controlled by the cluster.

System Services

Since there is more than one type of system service (`systemd` and `lsb`), Pacemaker supports a special `service` alias which intelligently figures out which one applies to a given cluster node.

This is particularly useful when the cluster contains a mix of `systemd` and `lsb`.

If the `service` standard is specified, Pacemaker will try to find the named service as an LSB init script, and if none exists, a `systemd` unit file.

STONITH

The `stonith` standard is used for managing fencing devices, discussed later in *Fencing*.

2.5.2 Resource Properties

These values tell the cluster which resource agent to use for the resource, where to find that resource agent and what standards it conforms to.

Table 9: Properties of a Primitive Resource

Field	Description
<code>id</code>	Your name for the resource
<code>class</code>	The standard the resource agent conforms to. Allowed values: <code>lsb</code> , <code>ocf</code> , <code>service</code> , <code>stonith</code> , and <code>systemd</code>
<code>description</code>	Arbitrary text for user's use (ignored by Pacemaker)
<code>type</code>	The name of the Resource Agent you wish to use. E.g. <code>IPaddr</code> or <code>Filesystem</code>
<code>provider</code>	The OCF spec allows multiple vendors to supply the same resource agent. To use the OCF resource agents supplied by the Heartbeat project, you would specify <code>heartbeat</code> here.

The XML definition of a resource can be queried with the `crm_resource` tool. For example:

```
# crm_resource --resource Email --query-xml
```

might produce:

A system resource definition

```
<primitive id="Email" class="service" type="exim"/>
```

Note: One of the main drawbacks to system services (`lsb` and `systemd`) is that they do not allow parameters

An OCF resource definition

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <instance_attributes id="Public-IP-params">
    <nvpair id="Public-IP-ip" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

2.5.3 Resource Options

Resources have two types of options: *meta-attributes* and *instance attributes*. Meta-attributes apply to any type of resource, while instance attributes are specific to each resource agent.

Resource Meta-Attributes

Meta-attributes are used by the cluster to decide how a resource should behave and can be easily set using the `--meta` option of the `crm_resource` command.

Table 10: Meta-attributes of a Primitive Resource

Name	Type	Default	Description
priority	<i>score</i>	0	If not all resources can be active, the cluster will stop lower-priority resources in order to keep higher-priority ones active.
critical	<i>boolean</i>	true	Use this value as the default for influence in all <i>colocation constraints</i> involving this resource, as well as in the implicit colocation constraints created if this resource is in a <i>group</i> . For details, see <i>Colocation Influence</i> . (since 2.1.0)
target-role	<i>enumeration</i>	Started	What state should the cluster attempt to keep this resource in? Allowed values: <ul style="list-style-type: none"> • Stopped: Force the resource to be stopped • Started: Allow the resource to be started (and in the case of <i>promotable</i> clone resources, promoted if appropriate) • Unpromoted: Allow the resource to be started, but only in the unpromoted role if the resource is <i>promotable</i> • Promoted: Equivalent to Started
is-managed	<i>boolean</i>	true	If false, the cluster will not start, stop, promote, or demote the resource on any node. Recurring actions for the resource are unaffected. Maintenance mode overrides this setting.

Continued on next page

Table 10 – continued from previous page

Name	Type	Default	Description
maintenance	<i>boolean</i>	false	If true, the cluster will not start, stop, promote, or demote the resource on any node, and will pause any recurring monitors (except those specifying <code>role</code> as <code>Stopped</code>). If true, the <i>maintenance-mode</i> cluster option or <i>maintenance</i> node attribute overrides this.
resource-stickiness	<i>score</i>	1 for individual clone instances, 0 for all other resources	A score that will be added to the current node when a resource is already active. This allows running resources to stay where they are, even if they would be placed elsewhere if they were being started from a stopped state.
requires	<i>enumeration</i>	<code>quorum</code> for resources with a <code>class</code> of <code>stonith</code> , otherwise <code>unfencing</code> if <code>unfencing</code> is active in the cluster, otherwise <code>fencing</code> if <code>stonith-enabled</code> is true, otherwise <code>quorum</code>	Conditions under which the resource can be started. Allowed values: <ul style="list-style-type: none"> • nothing: The cluster can always start this resource. • quorum: The cluster can start this resource only if a majority of the configured nodes are active. • fencing: The cluster can start this resource only if a majority of the configured nodes are active <i>and</i> any failed or unknown nodes have been <i>fenced</i>. • unfencing: The cluster can only start this resource if a majority of the configured nodes are active <i>and</i> any failed or unknown nodes have been <i>fenced</i> <i>and</i> only on nodes that have been <i>unfenced</i>.
migration-threshold	<i>score</i>	INFINITY	How many failures may occur for this resource on a node, before this node is marked ineligible to host this resource. A value of 0 indicates that this feature is disabled (the node will never be marked ineligible); by contrast, the cluster treats INFINITY (the default) as a very large but finite number. This option has an effect only if the failed operation specifies <code>on-fail</code> as <code>restart</code> (the default), and additionally for failed <code>start</code> operations, if the cluster property <code>start-failure-is-fatal</code> is <code>false</code> .

Continued on next page

Table 10 – continued from previous page

Name	Type	Default	Description
failure-timeout	<i>duration</i>	0	Ignore previously failed resource actions after this much time has passed without new failures (potentially allowing the resource back to the node on which it failed, if it previously reached its <code>migration-threshold</code> there). A value of 0 indicates that failures do not expire. WARNING: If this value is low, and pending cluster activity prevents the cluster from responding to a failure within that time, then the failure will be ignored completely and will not cause recovery of the resource, even if a recurring action continues to report failure. It should be at least greater than the longest <i>action timeout</i> for all resources in the cluster. A value in hours or days is reasonable.
multiple-active	<i>enumeration</i>	stop_start	What should the cluster do if it ever finds the resource active on more than one node? Allowed values: <ul style="list-style-type: none"> • <code>block</code>: mark the resource as unmanaged • <code>stop_only</code>: stop all active instances and leave them that way • <code>stop_start</code>: stop all active instances and start the resource in one location only • <code>stop_unexpected</code>: stop all active instances except where the resource should be active (this should be used only when extra instances are not expected to disrupt existing instances, and the resource agent’s monitor of an existing instance is capable of detecting any problems that could be caused; note that any resources ordered after this will still need to be restarted) (<i>since 2.1.3</i>)
allow-migrate	<i>boolean</i>	true for <code>ocf:pacemaker:remote</code> resources, false otherwise	Whether the cluster should try to “live migrate” this resource when it needs to be moved (see <i>Migrating Resources</i>)
allow-unhealthy-nodes	<i>boolean</i>	false	Whether the resource should be able to run on a node even if the node’s health score would otherwise prevent it (see <i>Tracking Node Health</i>) (<i>since 2.1.3</i>)
container-attribute-target	<i>enumeration</i>		Specific to bundle resources; see <i>Bundle Node Attributes</i>

As an example of setting resource options, if you performed the following commands on an LSB Email

resource:

```
# crm_resource --meta --resource Email --set-parameter priority --parameter-value 100
# crm_resource -m -r Email -p multiple-active -v block
```

the resulting resource definition might be:

An LSB resource with cluster options

```
<primitive id="Email" class="lsb" type="exim">
  <meta_attributes id="Email-meta_attributes">
    <nvpair id="Email-meta_attributes-priority" name="priority" value="100"/>
    <nvpair id="Email-meta_attributes-multiple-active" name="multiple-active" value="block"/>
  </meta_attributes>
</primitive>
```

In addition to the cluster-defined meta-attributes described above, you may also configure arbitrary meta-attributes of your own choosing. Most commonly, this would be done for use in *rules*. For example, an IT department might define a custom meta-attribute to indicate which company department each resource is intended for. To reduce the chance of name collisions with cluster-defined meta-attributes added in the future, it is recommended to use a unique, organization-specific prefix for such attributes.

Setting Global Defaults for Resource Meta-Attributes

To set a default value for a resource option, add it to the `rsc_defaults` section with `crm_attribute`. For example,

```
# crm_attribute --type rsc_defaults --name is-managed --update false
```

would prevent the cluster from starting or stopping any of the resources in the configuration (unless of course the individual resources were specifically enabled by having their `is-managed` set to `true`).

Resource Instance Attributes

The resource agents of some resource standards (lsb and *systemd* *not* among them) can be given parameters which determine how they behave and which instance of a service they control.

If your resource agent supports parameters, you can add them with the `crm_resource` command. For example,

```
# crm_resource --resource Public-IP --set-parameter ip --parameter-value 192.0.2.2
```

would create an entry in the resource like this:

An example OCF resource with instance attributes

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```


For an OCF resource, the result would be an environment variable called `OCF_RESKEY_ip` with a value of `192.0.2.2`.

The list of instance attributes supported by an OCF resource agent can be found by calling the resource agent with the `meta-data` command. The output contains an XML description of all the supported attributes, their purpose and default values.

Displaying the metadata for the Dummy resource agent template

```
# export OCF_ROOT=/usr/lib/ocf
# $OCF_ROOT/resource.d/pacemaker/Dummy meta-data
```

```
<?xml version="1.0"?>
<!DOCTYPE resource-agent SYSTEM "ra-api-1.dtd">
<resource-agent name="Dummy" version="2.0">
<version>1.1</version>

<longdesc lang="en">
This is a dummy OCF resource agent. It does absolutely nothing except keep track
of whether it is running or not, and can be configured so that actions fail or
take a long time. Its purpose is primarily for testing, and to serve as a
template for resource agent writers.
</longdesc>
<shortdesc lang="en">Example stateless resource agent</shortdesc>

<parameters>
<parameter name="state" unique-group="state">
<longdesc lang="en">
Location to store the resource state in.
</longdesc>
<shortdesc lang="en">State file</shortdesc>
<content type="string" default="/var/run/Dummy-RESOURCE_ID.state" />
</parameter>

<parameter name="passwd" reloadable="1">
<longdesc lang="en">
Fake password field
</longdesc>
<shortdesc lang="en">Password</shortdesc>
<content type="string" default="" />
</parameter>

<parameter name="fake" reloadable="1">
<longdesc lang="en">
Fake attribute that can be changed to cause a reload
</longdesc>
<shortdesc lang="en">Fake attribute that can be changed to cause a reload</shortdesc>
<content type="string" default="dummy" />
</parameter>

<parameter name="op_sleep" reloadable="1">
<longdesc lang="en">
Number of seconds to sleep during operations. This can be used to test how
the cluster reacts to operation timeouts.
</longdesc>
<shortdesc lang="en">Operation sleep duration in seconds.</shortdesc>
<content type="string" default="0" />
</parameter>

<parameter name="fail_start_on" reloadable="1">
<longdesc lang="en">
```

Start, migrate from, and reload-agent actions will return failure if running on the host specified here, but the resource will run successfully anyway (future monitor calls will find it running). This can be used to test on-fail=ignore.

2.5.4 Pacemaker Remote Resources

Pacemaker Remote nodes are defined by resources.

Remote nodes

A remote node is defined by a connection resource using the special, built-in `ocf:pacemaker:remote` resource agent.

Table 11: `ocf:pacemaker:remote` Instance Attributes

Name	Type	Default	Description
server	<i>text</i>	resource ID	Hostname or IP address used to connect to the remote node. The remote executor on the remote node must be configured to accept connections on this address.
port	<i>port</i>	3121	TCP port on the remote node used for its Pacemaker Remote connection. The remote executor on the remote node must be configured to listen on this port.
reconnect_interval	<i>duration</i>	0	If positive, the cluster will attempt to reconnect to a remote node at this interval after an active connection has been lost. Otherwise, the cluster will attempt to reconnect immediately (after any fencing, if needed).

Guest Nodes

When configuring a virtual machine as a guest node, the virtual machine is created using one of the usual resource agents for that purpose (for example, `ocf:heartbeat:VirtualDomain` or `ocf:heartbeat:Xen`), with additional meta-attributes.

No restrictions are enforced on what agents may be used to create a guest node, but obviously the agent must create a distinct environment capable of running the remote executor and cluster resources. An additional requirement is that fencing the node hosting the guest node resource must be sufficient for ensuring the guest node is stopped. This means that not all hypervisors supported by `VirtualDomain` may be used to create guest nodes; if the guest can survive the hypervisor being fenced, it is unsuitable for use as a guest node.

Table 12: `Guest node meta-attributes`

Name	Type	Default	Description
remote-node	<i>text</i>		If specified, this resource defines a guest node using this node name. The guest must be configured to run the remote executor when it is started. This value <i>must not</i> be the same as any resource or node ID.

Continued on next page

Table 12 – continued from previous page

Name	Type	Default	Description
remote-addr	<i>text</i>	value of <code>remote-node</code>	If <code>remote-node</code> is specified, the hostname or IP address used to connect to the guest. The remote executor on the guest must be configured to accept connections on this address.
remote-port	<i>port</i>	3121	If <code>remote-node</code> is specified, the port on the guest used for its Pacemaker Remote connection. The remote executor on the guest must be configured to listen on this port.
remote-connect-timeout	<i>timeout</i>	60s	If <code>remote-node</code> is specified, how long before a pending guest connection will time out.
remote-allow-migrate	<i>boolean</i>	true	If <code>remote-node</code> is specified, this acts as the <code>allow-migrate</code> meta-attribute for its implicitly created remote connection resource (<code>ocf:pacemaker:remote</code>).

Removing Pacemaker Remote Nodes

If the resource creating a remote node connection or guest node is removed from the configuration, status output may continue to show the affected node (as offline).

If you want to get rid of that output, run the following command, replacing `$NODE_NAME` appropriately:

```
# crm_node --force --remove $NODE_NAME
```

Warning: Be absolutely sure that there are no references to the node's resource in the configuration before running the above command.

2.6 Resource Operations

Operations are actions the cluster can perform on a resource by calling the resource agent. Resource agents must support certain common operations such as start, stop, and monitor, and may implement any others.

Operations may be explicitly configured for two purposes: to override defaults for options (such as timeout) that the cluster will use whenever it initiates the operation, and to run an operation on a recurring basis (for example, to monitor the resource for failure).

An OCF resource with a non-default start timeout

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <operations>
    <op id="Public-IP-start" name="start" timeout="60s"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

Pacemaker identifies operations by a combination of name and interval, so this combination must be unique for each resource. That is, you should not configure two operations for the same resource with the same name and interval.

2.6.1 Operation Properties

The `id`, `name`, `interval`, and `role` operation properties may be specified only as XML attributes of the `op` element. Other operation properties may be specified in any of the following ways, from highest precedence to lowest:

- directly in the `op` element as an XML attribute
- in an `nvpair` element within a `meta_attributes` element within the `op` element
- in an `nvpair` element within a `meta_attributes` element within *operation defaults*

If not specified, the default from the table below is used.

Table 13: Operation Properties

Name	Type	Default	Description
<code>id</code>	<i>id</i>		A unique identifier for the XML element (<i>required</i>)
<code>name</code>	<i>text</i>		An action name supported by the resource agent (<i>required</i>)
<code>interval</code>	<i>duration</i>	0	If this is a positive value, Pacemaker will schedule recurring instances of this operation at the given interval (which makes sense only with <i>name</i> set to <i>monitor</i>). If this is 0, Pacemaker will apply other properties configured for this operation to instances that are scheduled as needed during normal cluster operation. (<i>required</i>)
<code>description</code>	<i>text</i>		Arbitrary text for user's use (ignored by Pacemaker)
<code>role</code>	<i>enumeration</i>		If this is set, the operation configuration applies only on nodes where the cluster expects the resource to be in the specified role. This makes sense only for recurring monitors. Allowed values: Started , Stopped , and in the case of <i>promotable clone resources</i> , Unpromoted and Promoted .
<code>timeout</code>	<i>timeout</i>	20s	If resource agent execution does not complete within this amount of time, the action will be considered failed. Note: timeouts for fencing agents are handled specially (see the <i>Fencing</i> chapter).

Continued on next page

Table 13 – continued from previous page

Name	Type	Default	Description
on-fail	<i>enumeration</i>	<ul style="list-style-type: none"> • If <code>name</code> is <code>stop: fence</code> if <code>stonith-enabled</code> is true, otherwise <code>block</code> • If <code>name</code> is <code>demote: on-fail</code> of the <code>monitor</code> action with <code>role</code> set to <code>Promoted</code>, if present, enabled, and configured to a value other than <code>demote</code>, or <code>restart</code> otherwise • Otherwise: <code>restart</code> 	<p>How the cluster should respond to a failure of this action. Allowed values:</p> <ul style="list-style-type: none"> • <code>ignore</code>: Pretend the resource did not fail • <code>block</code>: Do not perform any further operations on the resource • <code>stop</code>: Stop the resource and leave it stopped • <code>demote</code>: Demote the resource, without a full restart. This is valid only for <code>promote</code> actions, and for <code>monitor</code> actions with both a nonzero <code>interval</code> and <code>role</code> set to <code>Promoted</code>; for any other action, a configuration error will be logged, and the default behavior will be used. (<i>since 2.0.5</i>) • <code>restart</code>: Stop the resource, and start it again if allowed (possibly on a different node) • <code>fence</code>: Fence the node on which the resource failed • <code>standby</code>: Put the node on which the resource failed in standby mode (forcing <i>all</i> resources away)
enabled	<i>boolean</i>	true	<p>If <code>false</code>, ignore this operation definition. This does not suppress all actions of this type, but is typically used to pause a recurring monitor. This can complement the resource being unmanaged (<i>is-managed</i> set to <code>false</code>), which does not stop recurring operations. Maintenance mode, which does stop configured monitors, overrides this setting.</p>

Continued on next page

Table 13 – continued from previous page

Name	Type	Default	Description
interval-origin	<i>ISO 8601</i>		If set for a recurring action, the action will be scheduled for this time plus a multiple of the action's interval, rather than immediately after the resource gains the monitored role. For example, you might schedule an in-depth monitor to run once per day outside business hours, by setting this to the desired time (on any date) and setting interval to 24h. At most one of interval-origin and start-delay may be set.
start-delay	<i>duration</i>		If set, the cluster will wait this long before running the action (for the first time, if recurring). This is an advanced option that should generally be avoided. It can be useful for a recurring monitor if a resource agent incorrectly returns success from start before the service is actually ready, and the agent can't be corrected, or for a start action if a service takes a very long time to start, and you don't want to block the cluster from responding to other events during that time. If this delay is longer than 5 minutes, the cluster will pretend that the action succeeded when it is first scheduled for the purpose of other actions needed, then act on the result when it actually runs. At most one of interval-origin and start-delay may be set.
record-pending	<i>boolean</i>	true	Operation results are always recorded when the operation completes (successful or not). If this is true , operations will also be recorded when initiated, so that status output can indicate that the operation is in progress. (<i>deprecated since 3.0.0</i>)

Note: Only one action can be configured for any given combination of **name** and **interval**.

Note: When **on-fail** is set to **demote**, recovery from failure by a successful demote causes the cluster to recalculate whether and where a new instance should be promoted. The node with the failure is eligible, so if promotion scores have not changed, it will be promoted again.

There is no direct equivalent of **migration-threshold** for the promoted role, but the same effect can be

achieved with a location constraint using a *rule* with a node attribute expression for the resource's fail count. For example, to immediately ban the promoted role from a node with any failed promote or promoted instance monitor:

```
<rsc_location id="loc1" rsc="my_primitive">
  <rule id="rule1" score="-INFINITY" role="Promoted" boolean-op="or">
    <expression id="expr1" attribute="fail-count-my_primitive#promote_0"
      operation="gte" value="1"/>
    <expression id="expr2" attribute="fail-count-my_primitive#monitor_10000"
      operation="gte" value="1"/>
  </rule>
</rsc_location>
```

This example assumes that there is a promotable clone of the `my_primitive` resource (note that the primitive name, not the clone name, is used in the rule), and that there is a recurring 10-second-interval monitor configured for the promoted role (fail count attributes specify the interval in milliseconds).

2.6.2 Monitoring Resources for Failure

When Pacemaker first starts a resource, it runs one-time `monitor` operations (referred to as *probes*) to ensure the resource is running where it's supposed to be, and not running where it's not supposed to be. (This behavior can be affected by the `resource-discovery` location constraint property.)

Other than those initial probes, Pacemaker will *not* (by default) check that the resource continues to stay healthy¹. You must configure `monitor` operations explicitly to perform these checks.

An OCF resource with a recurring health check

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <operations>
    <op id="Public-IP-start" name="start" timeout="60s"/>
    <op id="Public-IP-monitor" name="monitor" interval="60s"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

By default, a `monitor` operation will ensure that the resource is running where it is supposed to. The `target-role` property can be used for further checking.

For example, if a resource has one `monitor` operation with `interval=10` `role=Started` and a second `monitor` operation with `interval=11` `role=Stopped`, the cluster will run the first monitor on any nodes it thinks *should* be running the resource, and the second monitor on any nodes that it thinks *should not* be running the resource (for the truly paranoid, who want to know when an administrator manually starts a service by mistake).

Note: Currently, monitors with `role=Stopped` are not implemented for *clone* resources.

¹ Currently, anyway. Automatic monitoring operations may be added in a future version of Pacemaker.

2.6.3 Custom Recurring Operations

Typically, only `monitor` operations should be configured as recurring. However, it is possible to implement a custom action name in an OCF agent and then configure that as a recurring operation.

This could be useful, for example, to run a report, rotate a log, or clean temporary files related to a particular service.

Failures of custom recurring operations will be ignored by the cluster and will not be reported in cluster status (*since 3.0.0; previously, they would be treated like failed monitors*). A fail count and last failure timestamp will be recorded as transient node attributes, and those node attributes will be erased by the `crm_resource --cleanup` command.

2.6.4 Setting Global Defaults for Operations

You can change the global default values for operation properties in a given cluster. These are defined in an `op_defaults` section of the CIB's configuration section, and can be set with `crm_attribute`. For example,

```
# crm_attribute --type op_defaults --name timeout --update 20s
```

would default each operation's `timeout` to 20 seconds. If an operation's definition also includes a value for `timeout`, then that value would be used for that operation instead.

2.6.5 When Implicit Operations Take a Long Time

The cluster will always perform a number of implicit operations: `start`, `stop` and a non-recurring `monitor` operation used at startup to check whether the resource is already active. If one of these is taking too long, then you can create an entry for them and specify a longer timeout.

An OCF resource with custom timeouts for its implicit actions

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <operations>
    <op id="public-ip-startup" name="monitor" interval="0" timeout="90s"/>
    <op id="public-ip-start" name="start" interval="0" timeout="180s"/>
    <op id="public-ip-stop" name="stop" interval="0" timeout="15min"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

2.6.6 Multiple Monitor Operations

Provided no two operations (for a single resource) have the same name and interval, you can have as many `monitor` operations as you like. In this way, you can do a superficial health check every minute and progressively more intense ones at higher intervals.

To tell the resource agent what kind of check to perform, you need to provide each monitor with a different value for a common parameter. The OCF standard creates a special parameter called `OCF_CHECK_LEVEL` for

this purpose and dictates that it is “made available to the resource agent without the normal OCF_RESKEY prefix”.

Whatever name you choose, you can specify it by adding an `instance_attributes` block to the `op` tag. It is up to each resource agent to look for the parameter and decide how to use it.

An OCF resource with two recurring health checks, performing different levels of checks specified via OCF_CHECK_LEVEL.

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <operations>
    <op id="public-ip-health-60" name="monitor" interval="60">
      <instance_attributes id="params-public-ip-depth-60">
        <nvpair id="public-ip-depth-60" name="OCF_CHECK_LEVEL" value="10"/>
      </instance_attributes>
    </op>
    <op id="public-ip-health-300" name="monitor" interval="300">
      <instance_attributes id="params-public-ip-depth-300">
        <nvpair id="public-ip-depth-300" name="OCF_CHECK_LEVEL" value="20"/>
      </instance_attributes>
    </op>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-level" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

2.6.7 Disabling a Monitor Operation

The easiest way to stop a recurring monitor is to just delete it. However, there can be times when you only want to disable it temporarily. In such cases, simply add `enabled=false` to the operation’s definition.

Example of an OCF resource with a disabled health check

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <operations>
    <op id="public-ip-check" name="monitor" interval="60s" enabled="false"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

This can be achieved from the command line by executing:

```
# cibadmin --modify --xml-text '<op id="public-ip-check" enabled="false"/>'
```

Once you’ve done whatever you needed to do, you can then re-enable it with

```
# cibadmin --modify --xml-text '<op id="public-ip-check" enabled="true"/>'
```

2.6.8 Handling Resource Failure

By default, Pacemaker will attempt to recover failed resources by restarting them. However, failure recovery is highly configurable.

Failure Counts

Pacemaker tracks resource failures for each combination of node, resource, and operation (start, stop, monitor, etc.).

You can query the fail count for a particular node, resource, and/or operation using the `crm_failcount` command. For example, to see how many times the 10-second monitor for `myrsc` has failed on `node1`, run:

```
# crm_failcount --query -r myrsc -N node1 -n monitor -I 10s
```

If you omit the node, `crm_failcount` will use the local node. If you omit the operation and interval, `crm_failcount` will display the sum of the fail counts for all operations on the resource.

You can use `crm_resource --cleanup` or `crm_failcount --delete` to clear fail counts. For example, to clear the above monitor failures, run:

```
# crm_resource --cleanup -r myrsc -N node1 -n monitor -I 10s
```

If you omit the resource, `crm_resource --cleanup` will clear failures for all resources. If you omit the node, it will clear failures on all nodes. If you omit the operation and interval, it will clear the failures for all operations on the resource.

Note: Even when cleaning up only a single operation, all failed operations will disappear from the status display. This allows us to trigger a re-check of the resource's current status.

Higher-level tools may provide other commands for querying and clearing fail counts.

The `crm_mon` tool shows the current cluster status, including any failed operations. To see the current fail counts for any failed resources, call `crm_mon` with the `--failcounts` option. This shows the fail counts per resource (that is, the sum of any operation fail counts for the resource).

Failure Response

Normally, if a running resource fails, pacemaker will try to stop it and start it again. Pacemaker will choose the best location to start it each time, which may be the same node that it failed on.

However, if a resource fails repeatedly, it is possible that there is an underlying problem on that node, and you might desire trying a different node in such a case. Pacemaker allows you to set your preference via the `migration-threshold` resource meta-attribute.²

If you define `migration-threshold` to N for a resource, it will be banned from the original node after N failures there.

Note: The `migration-threshold` is per *resource*, even though fail counts are tracked per *operation*. The operation fail counts are added together to compare against the `migration-threshold`.

² The naming of this option was perhaps unfortunate as it is easily confused with live migration, the process of moving a resource from one node to another without stopping it. Xen virtual guests are the most common example of resources that can be migrated in this manner.

By default, fail counts remain until manually cleared by an administrator using `crm_resource --cleanup` or `crm_failcount --delete` (hopefully after first fixing the failure's cause). It is possible to have fail counts expire automatically by setting the `failure-timeout` resource meta-attribute.

Important: A successful operation does not clear past failures. If a recurring monitor operation fails once, succeeds many times, then fails again days later, its fail count is 2. Fail counts are cleared only by manual intervention or failure timeout.

For example, setting `migration-threshold` to 2 and `failure-timeout` to 60s would cause the resource to move to a new node after 2 failures, and allow it to move back (depending on stickiness and constraint scores) after one minute.

Note: `failure-timeout` is measured since the most recent failure. That is, older failures do not individually time out and lower the fail count. Instead, all failures are timed out simultaneously (and the fail count is reset to 0) if there is no new failure for the timeout period.

There are two exceptions to the migration threshold: when a resource either fails to start or fails to stop.

If the cluster property `start-failure-is-fatal` is set to `true` (which is the default), start failures cause the fail count to be set to `INFINITY` and thus always cause the resource to move immediately.

Stop failures are slightly different and crucial. If a resource fails to stop and fencing is enabled, then the cluster will fence the node in order to be able to start the resource elsewhere. If fencing is disabled, then the cluster has no way to continue and will not try to start the resource elsewhere, but will try to stop it again after any failure timeout or clearing.

2.6.9 Reloading an Agent After a Definition Change

The cluster automatically detects changes to the configuration of active resources. The cluster's normal response is to stop the service (using the old definition) and start it again (with the new definition). This works, but some resource agents are smarter and can be told to use a new set of options without restarting.

To take advantage of this capability, the resource agent must:

- Implement the `reload-agent` action. What it should do depends completely on your application!

Note: Resource agents may also implement a `reload` action to make the managed service reload its own *native* configuration. This is different from `reload-agent`, which makes effective changes in the resource's *Pacemaker* configuration (specifically, the values of the agent's reloadable parameters).

- Advertise the `reload-agent` operation in the `actions` section of its meta-data.
- Set the `reloadable` attribute to 1 in the `parameters` section of its meta-data for any parameters eligible to be reloaded after a change.

Once these requirements are satisfied, the cluster will automatically know to reload the resource (instead of restarting) when a reloadable parameter changes.

Note: Metadata will not be re-read unless the resource needs to be started. If you edit the agent of an already active resource to set a parameter reloadable, the resource may restart the first time the parameter value changes.

Note: If both a reloadable and non-reloadable parameter are changed simultaneously, the resource will be restarted.

2.6.10 Migrating Resources

Normally, when the cluster needs to move a resource, it fully restarts the resource (that is, it stops the resource on the current node and starts it on the new node).

However, some types of resources, such as many virtual machines, are able to move to another location without loss of state (often referred to as live migration or hot migration). In pacemaker, this is called live migration. Pacemaker can be configured to migrate a resource when moving it, rather than restarting it.

Not all resources are able to migrate; see the *migration checklist* below. Even those that can, won't do so in all situations. Conceptually, there are two requirements from which the other prerequisites follow:

- The resource must be active and healthy at the old location; and
- everything required for the resource to run must be available on both the old and new locations.

The cluster is able to accommodate both *push* and *pull* migration models by requiring the resource agent to support two special actions: `migrate_to` (performed on the current location) and `migrate_from` (performed on the destination).

In push migration, the process on the current location transfers the resource to the new location where it is later activated. In this scenario, most of the work would be done in the `migrate_to` action and, if anything, the activation would occur during `migrate_from`.

Conversely for pull, the `migrate_to` action is practically empty and `migrate_from` does most of the work, extracting the relevant resource state from the old location and activating it.

There is no wrong or right way for a resource agent to implement migration, as long as it works.

Migration Checklist

- The resource may not be a clone.
- The resource agent standard must be OCF.
- The resource must not be in a failed or degraded state.
- The resource agent must support `migrate_to` and `migrate_from` actions, and advertise them in its meta-data.
- The resource must have the `allow-migrate` meta-attribute set to `true` (which is not the default).

If an otherwise migratable resource depends on another resource via an ordering constraint, there are special situations in which it will be restarted rather than migrated.

For example, if the resource depends on a clone, and at the time the resource needs to be moved, the clone has instances that are stopping and instances that are starting, then the resource will be restarted. The scheduler is not yet able to model this situation correctly and so takes the safer (if less optimal) path.

Also, if a migratable resource depends on a non-migratable resource, and both need to be moved, the migratable resource will be restarted.

2.7 Resource Constraints

2.7.1 Deciding Which Nodes a Resource Can Run On

Location constraints tell the cluster which nodes a resource can run on.

There are two alternative strategies. One way is to say that, by default, resources can run anywhere, and then the location constraints specify nodes that are not allowed (an *opt-out* cluster). The other way is to start with nothing able to run anywhere, and use location constraints to selectively enable allowed nodes (an *opt-in* cluster).

Whether you should choose opt-in or opt-out depends on your personal preference and the make-up of your cluster. If most of your resources can run on most of the nodes, then an opt-out arrangement is likely to result in a simpler configuration. On the other-hand, if most resources can only run on a small subset of nodes, an opt-in configuration might be simpler.

Location Properties

Table 14: Attributes of a `rsc_location` Element

Name	Type	Default	Description
<code>id</code>	<i>id</i>		A unique name for the constraint (required)
<code>rsc</code>	<i>id</i>		The name of the resource to which this constraint applies. A location constraint must either have a <code>rsc</code> , have a <code>rsc-pattern</code> , or contain at least one resource set.
<code>rsc-pattern</code>	<i>text</i>		A pattern matching the names of resources to which this constraint applies. The syntax is the same as POSIX extended regular expressions, with the addition of an initial <code>!</code> indicating that resources <i>not</i> matching the pattern are selected. If the regular expression contains submatches, and the constraint contains a <i>rule</i> , the submatches can be referenced as <code>%1</code> through <code>%9</code> in the rule's <code>score-attribute</code> or a rule expression's <code>attribute</code> (see <i>Specifying location scores using pattern submatches</i>). A location constraint must either have a <code>rsc</code> , have a <code>rsc-pattern</code> , or contain at least one resource set.
<code>node</code>	<i>text</i>		The name of the node to which this constraint applies. A location constraint must either have a <code>node</code> and <code>score</code> , or contain at least one rule.
<code>score</code>	<i>score</i>		Positive values indicate a preference for running the affected resource(s) on <code>node</code> – the higher the value, the stronger the preference. Negative values indicate the resource(s) should avoid this node (a value of <code>-INFINITY</code> changes “should” to “must”). A location constraint must either have a <code>node</code> and <code>score</code> , or contain at least one rule.

Continued on next page

Table 14 – continued from previous page

Name	Type	Default	Description
role	<i>enumeration</i>	Started	<p>This is significant only for <i>promotable clones</i>, is allowed only if <code>rsc</code> or <code>rsc-pattern</code> is set, and is ignored if the constraint contains a rule. Allowed values:</p> <ul style="list-style-type: none"> • Started or Unpromoted: The constraint affects the location of all instances of the resource. (A promoted instance must start in the unpromoted role before being promoted, so any location requirement for unpromoted instances also affects promoted instances.) • Promoted: The constraint does not affect the location of instances, but instead affects which of the instances will be promoted.
resource-discovery	<i>enumeration</i>	always	<p>Whether Pacemaker should perform resource discovery (that is, check whether the resource is already running) for this resource on this node. This should normally be left as the default, so that rogue instances of a service can be stopped when they are running where they are not supposed to be. However, there are two situations where disabling resource discovery is a good idea: when a service is not installed on a node, discovery might return an error (properly written OCF agents will not, so this is usually only seen with other agent types); and when Pacemaker Remote is used to scale a cluster to hundreds of nodes, limiting resource discovery to allowed nodes can significantly boost performance. Allowed values:</p> <ul style="list-style-type: none"> • always: Always perform resource discovery for the specified resource on this node. • never: Never perform resource discovery for the specified resource on this node. This option should generally be used with a <code>-INFINITY</code> score, although that is not strictly required. • exclusive: Perform resource discovery for the specified resource only on this node (and other nodes similarly marked as exclusive). Multiple location constraints using exclusive discovery for the same resource across different nodes creates a subset of nodes resource-discovery is exclusive to. If a resource is marked for exclusive discovery on one or more nodes, that resource is only allowed to be placed within that subset of nodes.

Warning: Setting `resource-discovery` to `never` or `exclusive` removes Pacemaker’s ability to detect and stop unwanted instances of a service running where it’s not supposed to be. It is up to the system administrator (you!) to make sure that the service can *never* be active on nodes without `resource-discovery` (such as by leaving the relevant software uninstalled).

Asymmetrical “Opt-In” Clusters

To create an opt-in cluster, start by preventing resources from running anywhere by default:

```
# crm_attribute --name symmetric-cluster --update false
```

Then start enabling nodes. The following fragment says that the web server prefers **sles-1**, the database prefers **sles-2** and both can fail over to **sles-3** if their most preferred node fails.

Opt-in location constraints for two resources

```
<constraints>
  <rsc_location id="loc-1" rsc="Webserver" node="sles-1" score="200"/>
  <rsc_location id="loc-2" rsc="Webserver" node="sles-3" score="0"/>
  <rsc_location id="loc-3" rsc="Database" node="sles-2" score="200"/>
  <rsc_location id="loc-4" rsc="Database" node="sles-3" score="0"/>
</constraints>
```

Symmetrical “Opt-Out” Clusters

To create an opt-out cluster, start by allowing resources to run anywhere by default:

```
# crm_attribute --name symmetric-cluster --update true
```

Then start disabling nodes. The following fragment is the equivalent of the above opt-in configuration.

Opt-out location constraints for two resources

```
<constraints>
  <rsc_location id="loc-1" rsc="Webserver" node="sles-1" score="200"/>
  <rsc_location id="loc-2-do-not-run" rsc="Webserver" node="sles-2" score="-INFINITY"/>
  <rsc_location id="loc-3-do-not-run" rsc="Database" node="sles-1" score="-INFINITY"/>
  <rsc_location id="loc-4" rsc="Database" node="sles-2" score="200"/>
</constraints>
```

What if Two Nodes Have the Same Score

If two nodes have the same score, then the cluster will choose one. This choice may seem random and may not be what was intended, however the cluster was not given enough information to know any better.

Constraints where a resource prefers two nodes equally

```
<constraints>
  <rsc_location id="loc-1" rsc="Webserver" node="sles-1" score="INFINITY"/>
  <rsc_location id="loc-2" rsc="Webserver" node="sles-2" score="INFINITY"/>
  <rsc_location id="loc-3" rsc="Database" node="sles-1" score="500"/>
  <rsc_location id="loc-4" rsc="Database" node="sles-2" score="300"/>
  <rsc_location id="loc-5" rsc="Database" node="sles-2" score="200"/>
</constraints>
```

In the example above, assuming no other constraints and an inactive cluster, **Webserver** would probably be placed on **sles-1** and **Database** on **sles-2**. It would likely have placed **Webserver** based on the node's

uname and **Database** based on the desire to spread the resource load evenly across the cluster. However other factors can also be involved in more complex configurations.

Specifying locations using pattern matching

A location constraint can affect all resources whose IDs match a given pattern. The following example bans resources named **ip-httdp**, **ip-asterisk**, **ip-gateway**, etc., from **node1**.

Location constraint banning all resources matching a pattern from one node

```
<constraints>
  <rsc_location id="ban-ips-from-node1" rsc-pattern="ip-.*" node="node1" score="-INFINITY"/>
</constraints>
```

2.7.2 Specifying the Order in which Resources Should Start/Stop

Ordering constraints tell the cluster the order in which certain resource actions should occur.

Important: Ordering constraints affect *only* the ordering of resource actions; they do *not* require that the resources be placed on the same node. If you want resources to be started on the same node *and* in a specific order, you need both an ordering constraint *and* a colocation constraint (see *Placing Resources Relative to other Resources*), or alternatively, a group (see *Groups - A Syntactic Shortcut*).

Ordering Properties

Table 15: Attributes of a `rsc_order` Element

Field	Default	Description
id		A unique name for the constraint
first		Name of the resource that the then resource depends on
then		Name of the dependent resource
first-action	start	The action that the first resource must complete before then-action can be initiated for the then resource. Allowed values: start , stop , promote , demote .
then-action	value of first-action	The action that the then resource can execute only after the first-action on the first resource has completed. Allowed values: start , stop , promote , demote .

Continued on next page

Table 15 – continued from previous page

Field	Default	Description
kind	Mandatory	How to enforce the constraint. Allowed values: <ul style="list-style-type: none"> • Mandatory: then-action will never be initiated for the then resource unless and until first-action successfully completes for the first resource. • Optional: The constraint applies only if both specified resource actions are scheduled in the same transition (that is, in response to the same cluster state). This means that then-action is allowed on the then resource regardless of the state of the first resource, but if both actions happen to be scheduled at the same time, they will be ordered. • Serialize: Ensure that the specified actions are never performed concurrently for the specified resources. First-action and then-action can be executed in either order, but one must complete before the other can be initiated. An example use case is when resource start-up puts a high load on the host.
symmetrical	TRUE for Mandatory and Optional kinds. FALSE for Serialize kind.	If true, the reverse of the constraint applies for the opposite action (for example, if B starts after A starts, then B stops before A stops). Serialize orders cannot be symmetrical.

Promote and demote apply to *promotable* clone resources.

Optional and mandatory ordering

Here is an example of ordering constraints where **Database** *must* start before **Webserver**, and **IP** *should* start before **Webserver** if they both need to be started:

Optional and mandatory ordering constraints

```
<constraints>
  <rsc_order id="order-1" first="IP" then="Webserver" kind="Optional"/>
  <rsc_order id="order-2" first="Database" then="Webserver" kind="Mandatory" />
</constraints>
```

Because the above example lets `symmetrical` default to TRUE, **Webserver** must be stopped before **Database** can be stopped, and **Webserver** should be stopped before **IP** if they both need to be stopped.

Symmetric and asymmetric ordering

A mandatory symmetric ordering of “start A then start B” implies not only that the start actions must be ordered, but that B is not allowed to be active unless A is active. For example, if the ordering is added to the configuration when A is stopped (due to target-role, failure, etc.) and B is already active, then B will be stopped.

By contrast, asymmetric ordering of “start A then start B” means the stops can occur in either order, which implies that B *can* remain active in the same situation.

2.7.3 Placing Resources Relative to other Resources

Colocation constraints tell the cluster that the location of one resource depends on the location of another one.

Colocation has an important side-effect: it affects the order in which resources are assigned to a node. Think about it: You can't place A relative to B unless you know where B is¹.

So when you are creating colocation constraints, it is important to consider whether you should colocate A with B, or B with A.

Important: Colocation constraints affect *only* the placement of resources; they do *not* require that the resources be started in a particular order. If you want resources to be started on the same node *and* in a specific order, you need both an ordering constraint (see *Specifying the Order in which Resources Should Start/Stop*) *and* a colocation constraint, or alternatively, a group (see *Groups - A Syntactic Shortcut*).

Colocation Properties

Table 16: Attributes of a `rsc_colocation` Constraint

Field	Default	Description
id		A unique name for the constraint (required).
rsc		The name of a resource that should be located relative to <code>with-rsc</code> . A colocation constraint must either contain at least one <i>resource set</i> , or specify both <code>rsc</code> and <code>with-rsc</code> .
with-rsc		The name of the resource used as the colocation target. The cluster will decide where to put this resource first and then decide where to put <code>rsc</code> . A colocation constraint must either contain at least one <i>resource set</i> , or specify both <code>rsc</code> and <code>with-rsc</code> .
node-attribute	#uname	If <code>rsc</code> and <code>with-rsc</code> are specified, this node attribute must be the same on the node running <code>rsc</code> and the node running <code>with-rsc</code> for the constraint to be satisfied. (For details, see <i>Colocation by Node Attribute</i> .)
score	0	Positive values indicate the resources should run on the same node. Negative values indicate the resources should run on different nodes. Values of +/- INFINITY change “should” to “must”.
rsc-role	Started	If <code>rsc</code> and <code>with-rsc</code> are specified, and <code>rsc</code> is a <i>promotable clone</i> , the constraint applies only to <code>rsc</code> instances in this role. Allowed values: <code>Started</code> , <code>Stopped</code> , <code>Promoted</code> , <code>Unpromoted</code> . For details, see <i>Promotable Clone Constraints</i> .
with-rsc-role	Started	If <code>rsc</code> and <code>with-rsc</code> are specified, and <code>with-rsc</code> is a <i>promotable clone</i> , the constraint applies only to <code>with-rsc</code> instances in this role. Allowed values: <code>Started</code> , <code>Stopped</code> , <code>Promoted</code> , <code>Unpromoted</code> . For details, see <i>Promotable Clone Constraints</i> .

Continued on next page

¹ While the human brain is sophisticated enough to read the constraint in any order and choose the correct one depending on the situation, the cluster is not quite so smart. Yet.

Table 16 – continued from previous page

Field	Default	Description
influence	value of critical meta-attribute for rsc	Whether to consider the location preferences of rsc when with-rsc is already active. Allowed values: true, false. For details, see <i>Colocation Influence</i> . (since 2.1.0)

Mandatory Placement

Mandatory placement occurs when the constraint’s score is **+INFINITY** or **-INFINITY**. In such cases, if the constraint can’t be satisfied, then the **rsc** resource is not permitted to run. For **score=INFINITY**, this includes cases where the **with-rsc** resource is not active.

If you need resource **A** to always run on the same machine as resource **B**, you would add the following constraint:

Mandatory colocation constraint for two resources

```
<rsc_colocation id="colocate" rsc="A" with-rsc="B" score="INFINITY"/>
```

Remember, because **INFINITY** was used, if **B** can’t run on any of the cluster nodes (for whatever reason) then **A** will not be allowed to run. Whether **A** is running or not has no effect on **B**.

Alternatively, you may want the opposite – that **A** *cannot* run on the same machine as **B**. In this case, use **score="-INFINITY"**.

Mandatory anti-colocation constraint for two resources

```
<rsc_colocation id="anti-colocate" rsc="A" with-rsc="B" score="-INFINITY"/>
```

Again, by specifying **-INFINITY**, the constraint is binding. So if the only place left to run is where **B** already is, then **A** may not run anywhere.

As with **INFINITY**, **B** can run even if **A** is stopped. However, in this case **A** also can run if **B** is stopped, because it still meets the constraint of **A** and **B** not running on the same node.

Advisory Placement

If mandatory placement is about “must” and “must not”, then advisory placement is the “I’d prefer if” alternative.

For colocation constraints with scores greater than **-INFINITY** and less than **INFINITY**, the cluster will try to accommodate your wishes, but may ignore them if other factors outweigh the colocation score. Those factors might include other constraints, resource stickiness, failure thresholds, whether other resources would be prevented from being active, etc.

Advisory colocation constraint for two resources

```
<rsc_colocation id="colocate-maybe" rsc="A" with-rsc="B" score="500"/>
```

Colocation by Node Attribute

The `node-attribute` property of a colocation constraint allows you to express the requirement, “these resources must be on similar nodes”.

As an example, imagine that you have two Storage Area Networks (SANs) that are not controlled by the cluster, and each node is connected to one or the other. You may have two resources `r1` and `r2` such that `r2` needs to use the same SAN as `r1`, but doesn't necessarily have to be on the same exact node. In such a case, you could define a *node attribute* named `san`, with the value `san1` or `san2` on each node as appropriate. Then, you could colocate `r2` with `r1` using `node-attribute` set to `san`.

Colocation Influence

By default, if A is colocated with B, the cluster will take into account A's preferences when deciding where to place B, to maximize the chance that both resources can run.

For a detailed look at exactly how this occurs, see [Colocation Explained](#).

However, if `influence` is set to `false` in the colocation constraint, this will happen only if B is inactive and needing to be started. If B is already active, A's preferences will have no effect on placing B.

An example of what effect this would have and when it would be desirable would be a nonessential reporting tool colocated with a resource-intensive service that takes a long time to start. If the reporting tool fails enough times to reach its migration threshold, by default the cluster will want to move both resources to another node if possible. Setting `influence` to `false` on the colocation constraint would mean that the reporting tool would be stopped in this situation instead, to avoid forcing the service to move.

The `critical` resource meta-attribute is a convenient way to specify the default for all colocation constraints and groups involving a particular resource.

Note: If a noncritical resource is a member of a group, all later members of the group will be treated as noncritical, even if they are marked as (or left to default to) critical.

2.7.4 Resource Sets

Resource sets allow multiple resources to be affected by a single constraint.

A set of 3 resources

```
<resource_set id="resource-set-example">
  <resource_ref id="A"/>
  <resource_ref id="B"/>
  <resource_ref id="C"/>
</resource_set>
```

Resource sets are valid inside `rsc_location`, `rsc_order` (see *Ordering Sets of Resources*), `rsc_colocation` (see *Colocating Sets of Resources*), and `rsc_ticket` (see *Configuring Ticket Dependencies*) constraints.

A resource set has a number of properties that can be set, though not all have an effect in all contexts.

Table 17: Attributes of a `resource_set` Element

Field	Default	Description
<code>id</code>		A unique name for the set (required)
<code>sequential</code>	<code>true</code>	Whether the members of the set must be acted on in order. Meaningful within <code>rsc_order</code> and <code>rsc_colocation</code> .
<code>require-all</code>	<code>true</code>	Whether all members of the set must be active before continuing. With the current implementation, the cluster may continue even if only one member of the set is started, but if more than one member of the set is starting at the same time, the cluster will still wait until all of those have started before continuing (this may change in future versions). Meaningful within <code>rsc_order</code> .
<code>role</code>		The constraint applies only to resource set members that are <i>Promotable clones</i> in this role. Meaningful within <code>rsc_location</code> , <code>rsc_colocation</code> and <code>rsc_ticket</code> . Allowed values: <code>Started</code> , <code>Promoted</code> , <code>Unpromoted</code> . For details, see <i>Promotable Clone Constraints</i> .
<code>action</code>	<code>start</code>	The action that applies to <i>all members</i> of the set. Meaningful within <code>rsc_order</code> . Allowed values: <code>start</code> , <code>stop</code> , <code>promote</code> , <code>demote</code> .
<code>score</code>		<i>Advanced use only.</i> Use a specific score for this set. Meaningful within <code>rsc_location</code> or <code>rsc_colocation</code> .
<code>kind</code>		<i>Advanced use only.</i> Use a specific kind for this set. Meaningful within <code>rsc_order</code> .

Anti-colocation Chains

Sometimes, you would like a set of resources to be anti-located with each other. For example, `resource1`, `resource2`, and `resource3` must all run on different nodes.

A straightforward approach would be to configure either separate colocations or a resource set, with `-INFINITY` scores between all the resources.

However, this will not work as expected.

Resource sets may in the future gain new syntax for this specific situation, but for now, a workaround is to use *utilization* instead of colocations to keep the resources apart. Create a utilization attribute for the anti-colocation, assign the same value to each resource, and give each node the capacity to run one resource.

2.7.5 Ordering Sets of Resources

A common situation is for an administrator to create a chain of ordered resources, such as:

A chain of ordered resources

```
<constraints>
  <rsc_order id="order-1" first="A" then="B" />
  <rsc_order id="order-2" first="B" then="C" />
  <rsc_order id="order-3" first="C" then="D" />
</constraints>
```

Visual representation of the four resources' start order for the above constraints



Ordered Set

To simplify this situation, *Resource Sets* can be used within ordering constraints:

A chain of ordered resources expressed as a set

```

<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-example" sequential="true">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_order>
</constraints>
  
```

While the set-based format is not less verbose, it is significantly easier to get right and maintain.

Important: If you use a higher-level tool, pay attention to how it exposes this functionality. Depending on the tool, creating a set **A B** may be equivalent to **A then B**, or **B then A**.

Ordering Multiple Sets

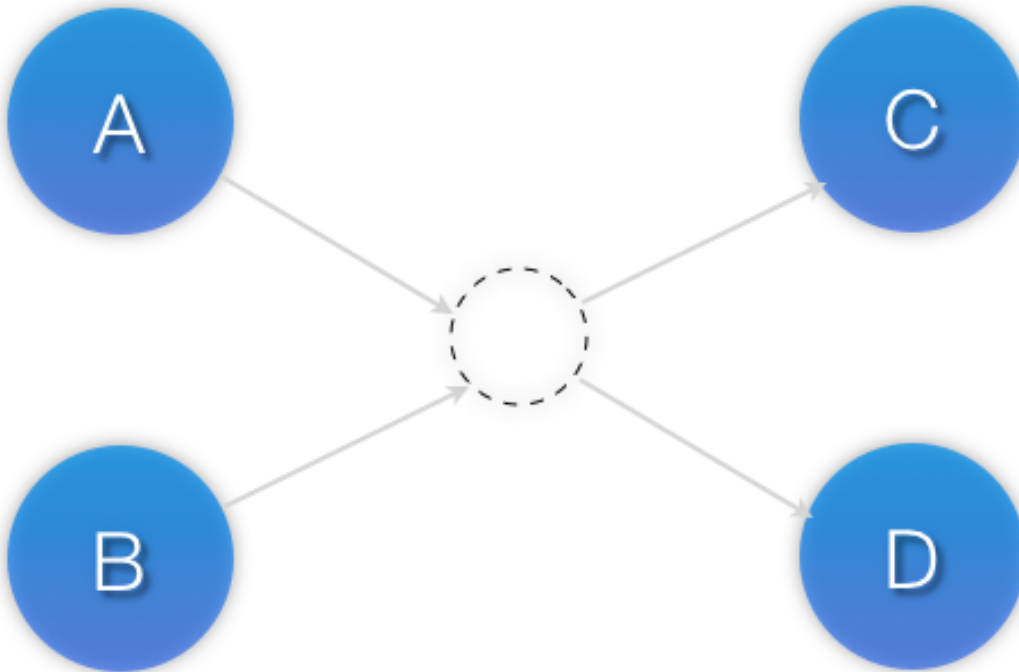
The syntax can be expanded to allow sets of resources to be ordered relative to each other, where the members of each individual set may be ordered or unordered (controlled by the `sequential` property). In the example below, **A** and **B** can both start in parallel, as can **C** and **D**, however **C** and **D** can only start once *both A and B* are active.

Ordered sets of unordered resources

```

<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-1" sequential="false">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="false">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_order>
</constraints>
  
```

Visual representation of the start order for two ordered sets of unordered resources

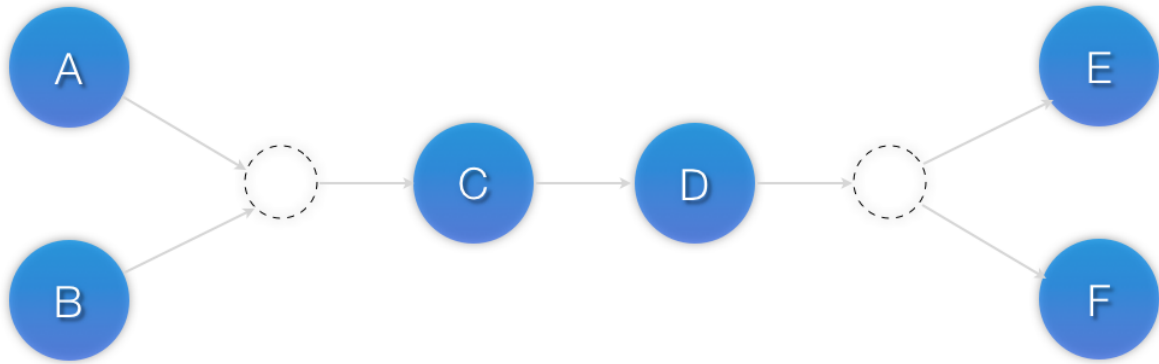


Of course either set – or both sets – of resources can also be internally ordered (by setting `sequential="true"`) and there is no limit to the number of sets that can be specified.

Advanced use of set ordering - Three ordered sets, two of which are internally unordered

```
<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-1" sequential="false">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="true">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
    <resource_set id="ordered-set-3" sequential="false">
      <resource_ref id="E"/>
      <resource_ref id="F"/>
    </resource_set>
  </rsc_order>
</constraints>
```

Visual representation of the start order for the three sets defined above



Important: An ordered set with `sequential=false` makes sense only if there is another set in the constraint. Otherwise, the constraint has no effect.

Resource Set OR Logic

The unordered set logic discussed so far has all been “AND” logic. To illustrate this take the 3 resource set figure in the previous section. Those sets can be expressed, **(A and B) then (C) then (D) then (E and F)**.

Say for example we want to change the first set, **(A and B)**, to use “OR” logic so the sets look like this: **(A or B) then (C) then (D) then (E and F)**. This functionality can be achieved through the use of the `require-all` option. This option defaults to TRUE which is why the “AND” logic is used by default. Setting `require-all=false` means only one resource in the set needs to be started before continuing on to the next set.

Resource Set “OR” logic: Three ordered sets, where the first set is internally unordered with “OR” logic

```

<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-1" sequential="false" require-all="false">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="true">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
    <resource_set id="ordered-set-3" sequential="false">
      <resource_ref id="E"/>
      <resource_ref id="F"/>
    </resource_set>
  </rsc_order>
</constraints>
  
```

Important: An ordered set with `require-all=false` makes sense only in conjunction with `sequential=false`. Think of it like this: `sequential=false` modifies the set to be an unordered set using “AND” logic by default, and adding `require-all=false` flips the unordered set’s “AND” logic to “OR” logic.

2.7.6 Colocating Sets of Resources

Another common situation is for an administrator to create a set of colocated resources.

The simplest way to do this is to define a resource group (see *Groups - A Syntactic Shortcut*), but that cannot always accurately express the desired relationships. For example, maybe the resources do not need to be ordered.

Another way would be to define each relationship as an individual constraint, but that causes a difficult-to-follow constraint explosion as the number of resources and combinations grow.

Colocation chain as individual constraints, where A is placed first, then B, then C, then D

```
<constraints>
  <rsc_colocation id="coloc-1" rsc="D" with-rsc="C" score="INFINITY"/>
  <rsc_colocation id="coloc-2" rsc="C" with-rsc="B" score="INFINITY"/>
  <rsc_colocation id="coloc-3" rsc="B" with-rsc="A" score="INFINITY"/>
</constraints>
```

To express complicated relationships with a simplified syntax², *resource sets* can be used within colocation constraints.

Equivalent colocation chain expressed using resource_set

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-example" sequential="true">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

Note: Within a `resource_set`, the resources are listed in the order they are *placed*, which is the reverse of the order in which they are *colocated*. In the above example, resource **A** is placed before resource **B**, which is the same as saying resource **B** is colocated with resource **A**.

As with individual constraints, a resource that can’t be active prevents any resource that must be colocated with it from being active. In both of the two previous examples, if **B** is unable to run, then both **C** and by inference **D** must remain stopped.

² which is not the same as saying easy to follow

Important: If you use a higher-level tool, pay attention to how it exposes this functionality. Depending on the tool, creating a set **A B** may be equivalent to **A with B**, or **B with A**.

Resource sets can also be used to tell the cluster that entire *sets* of resources must be colocated relative to each other, while the individual members within any one set may or may not be colocated relative to each other (determined by the set's `sequential` property).

In the following example, resources **B**, **C**, and **D** will each be colocated with **A** (which will be placed first). **A** must be able to run in order for any of the resources to run, but any of **B**, **C**, or **D** may be stopped without affecting any of the others.

Using colocated sets to specify a shared dependency

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-2" sequential="false">
      <resource_ref id="B"/>
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
    <resource_set id="colocated-set-1" sequential="true">
      <resource_ref id="A"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

Note: Pay close attention to the order in which resources and sets are listed. While the members of any one sequential set are placed first to last (i.e., the colocation dependency is last with first), multiple sets are placed last to first (i.e. the colocation dependency is first with last).

Important: A colocated set with `sequential="false"` makes sense only if there is another set in the constraint. Otherwise, the constraint has no effect.

There is no inherent limit to the number and size of the sets used. The only thing that matters is that in order for any member of one set in the constraint to be active, all members of sets listed after it must also be active (and naturally on the same node); and if a set has `sequential="true"`, then in order for one member of that set to be active, all members listed before it must also be active.

If desired, you can restrict the dependency to instances of promotable clone resources that are in a specific role, using the set's `role` property.

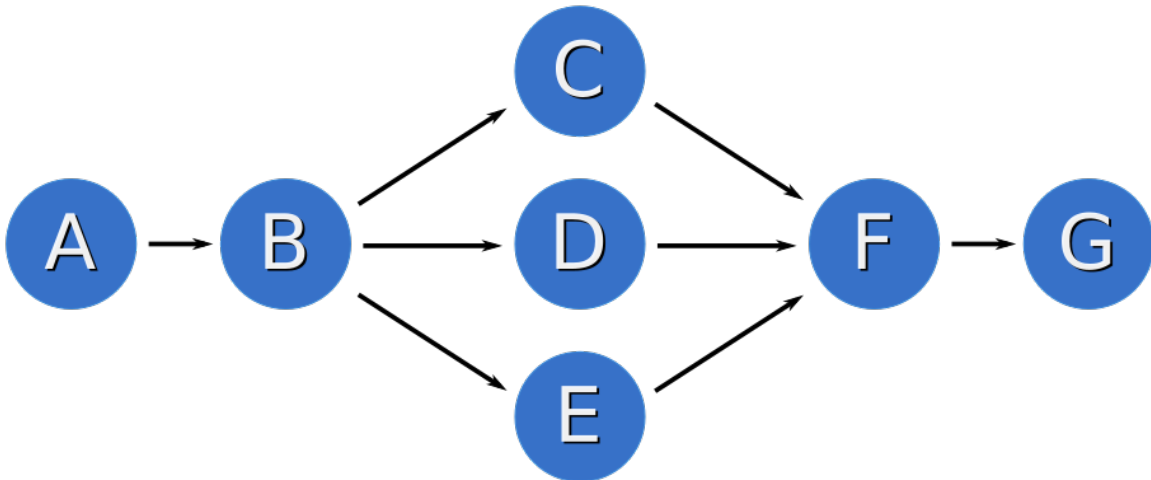
Colocation in which the members of the middle set have no interdependencies, and the last set listed applies only to promoted instances

```

<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-1" sequential="true">
      <resource_ref id="F"/>
      <resource_ref id="G"/>
    </resource_set>
    <resource_set id="colocated-set-2" sequential="false">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
      <resource_ref id="E"/>
    </resource_set>
    <resource_set id="colocated-set-3" sequential="true" role="Promoted">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
  </rsc_colocation>
</constraints>

```

Visual representation of the above example (resources are placed from left to right)



Note: Unlike ordered sets, colocated sets do not use the `require-all` option.

2.7.7 External Resource Dependencies

Sometimes, a resource will depend on services that are not managed by the cluster. An example might be a resource that requires a file system that is not managed by the cluster but mounted by `systemd` at boot time.

To accommodate this, the `pacemaker systemd` service depends on a normally empty target called `resource-agents-deps.target`. The system administrator may create a unit drop-in for that target specifying the dependencies, to ensure that the services are started before Pacemaker starts and stopped after Pacemaker stops.

Typically, this is accomplished by placing a unit file in the `/etc/systemd/system/resource-agents-deps.target.d` directory, with directives such as `Requires` and `After` specifying the dependencies as needed.

2.8 Fencing

2.8.1 What Is Fencing?

Fencing is the ability to make a node unable to run resources, even when that node is unresponsive to cluster commands.

Fencing is also known as *STONITH*, an acronym for “Shoot The Other Node In The Head”, since the most common fencing method is cutting power to the node. Another method is “fabric fencing”, cutting the node’s access to some capability required to run resources (such as network access or a shared disk).

2.8.2 Why Is Fencing Necessary?

Fencing protects your data from being corrupted by malfunctioning nodes or unintentional concurrent access to shared resources.

Fencing protects against the “split brain” failure scenario, where cluster nodes have lost the ability to reliably communicate with each other but are still able to run resources. If the cluster just assumed that uncommunicative nodes were down, then multiple instances of a resource could be started on different nodes.

The effect of split brain depends on the resource type. For example, an IP address brought up on two hosts on a network will cause packets to randomly be sent to one or the other host, rendering the IP useless. For a database or clustered file system, the effect could be much more severe, causing data corruption or divergence.

Fencing is also used when a resource cannot otherwise be stopped. If a resource fails to stop on a node, it cannot be started on a different node without risking the same type of conflict as split-brain. Fencing the original node ensures the resource can be safely started elsewhere.

Users may also configure the `on-fail` property of *Resource Operations* or the `loss-policy` property of *ticket constraints* to `fence`, in which case the cluster will fence the resource’s node if the operation fails or the ticket is lost.

2.8.3 Fence Devices

A *fence device* or *fencing device* is a special type of resource that provides the means to fence a node.

Examples of fencing devices include intelligent power switches and IPMI devices that accept SNMP commands to cut power to a node, and iSCSI controllers that allow SCSI reservations to be used to cut a node’s access to a shared disk.

Since fencing devices will be used to recover from loss of networking connectivity to other nodes, it is essential that they do not rely on the same network as the cluster itself, otherwise that network becomes a single point of failure.

Since loss of a node due to power outage is indistinguishable from loss of network connectivity to that node, it is also essential that at least one fence device for a node does not share power with that node. For example, an on-board IPMI controller that shares power with its host should not be used as the sole fencing device for that host.

Since fencing is used to isolate malfunctioning nodes, no fence device should rely on its target functioning properly. This includes, for example, devices that ssh into a node and issue a shutdown command (such devices might be suitable for testing, but never for production).

2.8.4 Fence Agents

A *fence agent* or *fencing agent* is a `stonith`-class resource agent.

The fence agent standard provides commands (such as `off` and `reboot`) that the cluster can use to fence nodes. As with other resource agent classes, this allows a layer of abstraction so that Pacemaker doesn't need any knowledge about specific fencing technologies – that knowledge is isolated in the agent.

Pacemaker supports two fence agent standards, both inherited from no-longer-active projects:

- Red Hat Cluster Suite (RHCS) style: These are typically installed in `/usr/sbin` with names starting with `fence_`.
- Linux-HA style: These typically have names starting with `external/`. Pacemaker can support these agents using the `fence_legacy` RHCS-style agent as a wrapper, *if* support was enabled when Pacemaker was built, which requires the `cluster-glue` library.

2.8.5 When a Fence Device Can Be Used

Fencing devices do not actually “run” like most services. Typically, they just provide an interface for sending commands to an external device.

Additionally, fencing may be initiated by Pacemaker, by other cluster-aware software such as DRBD or DLM, or manually by an administrator, at any point in the cluster life cycle, including before any resources have been started.

To accommodate this, Pacemaker does not require the fence device resource to be “started” in order to be used. Whether a fence device is started or not determines whether a node runs any recurring monitor for the device, and gives the node a slight preference for being chosen to execute fencing using that device.

By default, any node can execute any fencing device. If a fence device is disabled by setting its `target-role` to `Stopped`, then no node can use that device. If a location constraint with a negative score prevents a specific node from “running” a fence device, then that node will never be chosen to execute fencing using the device. A node may fence itself, but the cluster will choose that only if no other nodes can do the fencing.

A common configuration scenario is to have one fence device per target node. In such a case, users often configure anti-location constraints so that the target node does not monitor its own device.

2.8.6 Limitations of Fencing Resources

Fencing resources have certain limitations that other resource classes don't:

- They may have only one set of meta-attributes and one set of instance attributes.
- If *Rules* are used to determine fencing resource options, these might be evaluated only when first read, meaning that later changes to the rules will have no effect. Therefore, it is better to avoid confusion and not use rules at all with fencing resources.

These limitations could be revisited if there is sufficient user demand.

2.8.7 Special Meta-Attributes for Fencing Resources

The table below lists special resource meta-attributes that may be set for any fencing resource.

Table 18: Additional Properties of Fencing Resources

Field	Type	Default	Description
provides	string		Any special capability provided by the fence device. Currently, only one such capability is meaningful: <i>unfencing</i> .

2.8.8 Special Instance Attributes for Fencing Resources

The table below lists special instance attributes that may be set for any fencing resource (*not* meta-attributes, even though they are interpreted by Pacemaker rather than the fence agent). These are also listed in the man page for `pacemaker-fenced`.

Table 19: Additional Properties of Fencing Resources

Name	Type	Default	Description
stonith-timeout	<i>timeout</i>		This is not used by Pacemaker (see the <code>pcmk_reboot_timeout</code> , <code>pcmk_off_timeout</code> , etc., properties instead), but it may be used by Linux-HA fence agents.
pcmk_host_map	<i>text</i>		A mapping of node names to ports for devices that do not understand the node names. For example, <code>node1:1;node2:2,3</code> tells the cluster to use port 1 for <code>node1</code> and ports 2 and 3 for <code>node2</code> . If <code>pcmk_host_check</code> is explicitly set to <code>static-list</code> , either this or <code>pcmk_host_list</code> must be set. The port portion of the map may contain special characters such as spaces if preceded by a backslash (<i>since 2.1.2</i>).
pcmk_host_list	<i>text</i>		Comma-separated list of nodes that can be targeted by this device (for example, <code>node1,node2,node3</code>). If <code>pcmk_host_check</code> is <code>static-list</code> , either this or <code>pcmk_host_map</code> must be set.
pcmk_host_check	<i>text</i>	See <i>Default Check Type</i>	The method Pacemaker should use to determine which nodes can be targeted by this device. Allowed values: <ul style="list-style-type: none"> <code>static-list</code>: targets are listed in the <code>pcmk_host_list</code> or <code>pcmk_host_map</code> attribute <code>dynamic-list</code>: query the device via the agent's <code>list</code> action <code>status</code>: query the device via the agent's <code>status</code> action <code>none</code>: assume the device can fence any node

Continued on next page

Table 19 – continued from previous page

Name	Type	Default	Description
pcmk_delay_max	<i>duration</i>	0s	Enable a delay of no more than the time specified before executing fencing actions. Pacemaker derives the overall delay by taking the value of <code>pcmk_delay_base</code> and adding a random delay value such that the sum is kept below this maximum. This is sometimes used in two-node clusters to ensure that the nodes don't fence each other at the same time.
pcmk_delay_base	<i>text</i>	0s	Enable a static delay before executing fencing actions. This can be used, for example, in two-node clusters to ensure that the nodes don't fence each other, by having separate fencing resources with different values. The node that is fenced with the shorter delay will lose a fencing race. The overall delay introduced by pacemaker is derived from this value plus a random delay such that the sum is kept below the maximum delay. A single device can have different delays per node using a host map (<i>since 2.1.2</i>), for example <code>node1:0s;node2:5s</code> .
pcmk_action_limit	<i>integer</i>	1	The maximum number of actions that can be performed in parallel on this device. A value of -1 means unlimited. Node fencing actions initiated by the cluster (as opposed to an administrator running the <code>stonith_admin</code> tool or the fencer running recurring device monitors and <code>status</code> and <code>list</code> commands) are additionally subject to the <code>concurrent-fencing</code> cluster property.
pcmk_host_argument	<i>text</i>	<code>port</code> if the fence agent metadata advertises support for it, otherwise <code>plug</code> if supported, otherwise <code>none</code>	<i>Advanced use only.</i> Which parameter should be supplied to the fence agent to identify the node to be fenced. A value of <code>none</code> tells the cluster not to supply any additional parameters.
pcmk_reboot_action	<i>text</i>	<code>reboot</code>	<i>Advanced use only.</i> The command to send to the resource agent in order to reboot a node. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.
pcmk_reboot_timeout	<i>timeout</i>	60s	<i>Advanced use only.</i> Specify an alternate timeout (in seconds) to use for <code>reboot</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.

Continued on next page

Table 19 – continued from previous page

Name	Type	Default	Description
pcmk_reboot_retries	<i>integer</i>	2	<i>Advanced use only.</i> The maximum number of times to retry the reboot command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.
pcmk_off_action	<i>text</i>	off	<i>Advanced use only.</i> The command to send to the resource agent in order to shut down a node. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.
pcmk_off_timeout	<i>timeout</i>	60s	<i>Advanced use only.</i> Specify an alternate timeout (in seconds) to use for off actions instead of the value of stonith-timeout . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.
pcmk_off_retries	<i>integer</i>	2	<i>Advanced use only.</i> The maximum number of times to retry the off command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.
pcmk_list_action	<i>text</i>	list	<i>Advanced use only.</i> The command to send to the resource agent in order to list nodes. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.
pcmk_list_timeout	<i>timeout</i>	60s	<i>Advanced use only.</i> Specify an alternate timeout (in seconds) to use for list actions instead of the value of stonith-timeout . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.
pcmk_list_retries	<i>integer</i>	2	<i>Advanced use only.</i> The maximum number of times to retry the list command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.

Continued on next page

Table 19 – continued from previous page

Name	Type	Default	Description
<code>pcmk_monitor_action</code>	<i>text</i>	<code>monitor</code>	<i>Advanced use only.</i> The command to send to the resource agent in order to report extended status. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.
<code>pcmk_monitor_timeout</code>	<i>timeout</i>	60s	<i>Advanced use only.</i> Specify an alternate timeout (in seconds) to use for <code>monitor</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.
<code>pcmk_monitor_retries</code>	<i>integer</i>	2	<i>Advanced use only.</i> The maximum number of times to retry the <code>monitor</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.
<code>pcmk_status_action</code>	<i>text</i>	<code>status</code>	<i>Advanced use only.</i> The command to send to the resource agent in order to report status. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.
<code>pcmk_status_timeout</code>	<i>timeout</i>	60s	<i>Advanced use only.</i> Specify an alternate timeout (in seconds) to use for <code>status</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.
<code>pcmk_status_retries</code>	<i>integer</i>	2	<i>Advanced use only.</i> The maximum number of times to retry the <code>status</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.

2.8.9 Default Check Type

If the user does not explicitly configure `pcmk_host_check` for a fence device, a default value appropriate to other configured parameters will be used:

- If either `pcmk_host_list` or `pcmk_host_map` is configured, `static-list` will be used;
- otherwise, if the fence device supports the `list` action, and the first attempt at using `list` succeeds, `dynamic-list` will be used;

- otherwise, if the fence device supports the `status` action, `status` will be used;
- otherwise, `none` will be used.

2.8.10 Unfencing

With fabric fencing (such as cutting network or shared disk access rather than power), it is expected that the cluster will fence the node, and then a system administrator must manually investigate what went wrong, correct any issues found, then reboot (or restart the cluster services on) the node.

Once the node reboots and rejoins the cluster, some fabric fencing devices require an explicit command to restore the node's access. This capability is called *unfencing* and is typically implemented as the fence agent's `on` command.

If any cluster resource has `requires` set to `unfencing`, then that resource will not be probed or started on a node until that node has been unfenced.

2.8.11 Fencing and Quorum

In general, a cluster partition may execute fencing only if the partition has quorum, and the `stonith-enabled` cluster property is set to true. However, there are exceptions:

- The requirements apply only to fencing initiated by Pacemaker. If an administrator initiates fencing using the `stonith_admin` command, or an external application such as DLM initiates fencing using Pacemaker's C API, the requirements do not apply.
- A cluster partition without quorum is allowed to fence any active member of that partition. As a corollary, this allows a `no-quorum-policy` of `suicide` to work.
- If the `no-quorum-policy` cluster property is set to `ignore`, then quorum is not required to execute fencing of any node.

2.8.12 Fencing Timeouts

Fencing timeouts are complicated, since a single fencing operation can involve many steps, each of which may have a separate timeout.

Fencing may be initiated in one of several ways:

- An administrator may initiate fencing using the `stonith_admin` tool, which has a `--timeout` option (defaulting to 2 minutes) that will be used as the fence operation timeout.
- An external application such as DLM may initiate fencing using the Pacemaker C API. The application will specify the fence operation timeout in this case, which might or might not be configurable by the user.
- The cluster may initiate fencing itself. In this case, the `stonith-timeout` cluster property (defaulting to 1 minute) will be used as the fence operation timeout.

However fencing is initiated, the initiator contacts Pacemaker's fencer (`pacemaker-fenced`) to request fencing. This connection and request has its own timeout, separate from the fencing operation timeout, but usually happens very quickly.

The fencer will contact all fencers in the cluster to ask what devices they have available to fence the target node. The fence operation timeout will be used as the timeout for each of these queries.

Once a fencing device has been selected, the fencer will check whether any action-specific timeout has been configured for the device, to use instead of the fence operation timeout. For example, if `stonith-timeout`

is 60 seconds, but the fencing device has `pcmk_reboot_timeout` configured as 90 seconds, then a timeout of 90 seconds will be used for reboot actions using that device.

A device may have retries configured, in which case the timeout applies across all attempts. For example, if a device has `pcmk_reboot_retries` configured as 2, and the first reboot attempt fails, the second attempt will only have whatever time is remaining in the action timeout after subtracting how much time the first attempt used. This means that if the first attempt fails due to using the entire timeout, no further attempts will be made. There is currently no way to configure a per-attempt timeout.

If more than one device is required to fence a target, whether due to failure of the first device or a fencing topology with multiple devices configured for the target, each device will have its own separate action timeout.

For all of the above timeouts, the fencer will generally multiply the configured value by 1.2 to get an actual value to use, to account for time needed by the fencer's own processing.

Separate from the fencer's timeouts, some fence agents have internal timeouts for individual steps of their fencing process. These agents often have parameters to configure these timeouts, such as `login-timeout`, `shell-timeout`, or `power-timeout`. Many such agents also have a `disable-timeout` parameter to ignore their internal timeouts and just let Pacemaker handle the timeout. This causes a difference in retry behavior. If `disable-timeout` is not set, and the agent hits one of its internal timeouts, it will report that as a failure to Pacemaker, which can then retry. If `disable-timeout` is set, and Pacemaker hits a timeout for the agent, then there will be no time remaining, and no retry will be done.

2.8.13 Fence Devices Dependent on Other Resources

In some cases, a fence device may require some other cluster resource (such as an IP address) to be active in order to function properly.

This is obviously undesirable in general: fencing may be required when the depended-on resource is not active, or fencing may be required because the node running the depended-on resource is no longer responding.

However, this may be acceptable under certain conditions:

- The dependent fence device should not be able to target any node that is allowed to run the depended-on resource.
- The depended-on resource should not be disabled during production operation.
- The `concurrent-fencing` cluster property should be set to `true`. Otherwise, if both the node running the depended-on resource and some node targeted by the dependent fence device need to be fenced, the fencing of the node running the depended-on resource might be ordered first, making the second fencing impossible and blocking further recovery. With concurrent fencing, the dependent fence device might fail at first due to the depended-on resource being unavailable, but it will be retried and eventually succeed once the resource is brought back up.

Even under those conditions, there is one unlikely problem scenario. The DC always schedules fencing of itself after any other fencing needed, to avoid unnecessary repeated DC elections. If the dependent fence device targets the DC, and both the DC and a different node running the depended-on resource need to be fenced, the DC fencing will always fail and block further recovery. Note, however, that losing a DC node entirely causes some other node to become DC and schedule the fencing, so this is only a risk when a stop or other operation with `on-fail` set to `fencing` fails on the DC.

2.8.14 Configuring Fencing

Higher-level tools can provide simpler interfaces to this process, but using Pacemaker command-line tools, this is how you could configure a fence device.

1. Find the correct driver:

```
# stonith_admin --list-installed
```

Note: You may have to install packages to make fence agents available on your host. Searching your available packages for `fence-` is usually helpful. Ensure the packages providing the fence agents you require are installed on every cluster node.

2. Find the required parameters associated with the device (replacing `$AGENT_NAME` with the name obtained from the previous step):

```
# stonith_admin --metadata --agent $AGENT_NAME
```

3. Create a file called `stonith.xml` containing a primitive resource with a class of `stonith`, a type equal to the agent name obtained earlier, and a parameter for each of the values returned in the previous step.
4. If the device does not know how to fence nodes based on their `uname`, you may also need to set the special `pcmk_host_map` parameter. See *Special Instance Attributes for Fencing Resources* for details.
5. If the device does not support the `list` command, you may also need to set the special `pcmk_host_list` and/or `pcmk_host_check` parameters. See *Special Instance Attributes for Fencing Resources* for details.
6. If the device does not expect the target to be specified with the `port` parameter, you may also need to set the special `pcmk_host_argument` parameter. See *Special Instance Attributes for Fencing Resources* for details.
7. Upload it into the CIB using `cibadmin`:

```
# cibadmin --create --scope resources --xml-file stonith.xml
```

8. Set `stonith-enabled` to true:

```
# crm_attribute --type crm_config --name stonith-enabled --update true
```

9. Once the `stonith` resource is running, you can test it by executing the following, replacing `$NODE_NAME` with the name of the node to fence (although you might want to stop the cluster on that machine first):

```
# stonith_admin --reboot $NODE_NAME
```

Example Fencing Configuration

For this example, we assume we have a cluster node, `pcmk-1`, whose IPMI controller is reachable at the IP address `192.0.2.1`. The IPMI controller uses the username `testuser` and the password `abc123`.

1. Looking at what's installed, we may see a variety of available agents:

```
# stonith_admin --list-installed
```

```
(... some output omitted ...)  
fence_idrac  
fence_ilo3  
fence_ilo4  
fence_ilo5  
fence_imm
```

(continues on next page)

(continued from previous page)

```
fence_ipmilan
(... some output omitted ...)
```

Perhaps after some reading some man pages and doing some Internet searches, we might decide fence_ipmilan is our best choice.

2. Next, we would check what parameters fence_ipmilan provides:

```
# stonith_admin --metadata -a fence_ipmilan
```

```
<resource-agent name="fence_ipmilan" shortdesc="Fence agent for IPMI">
  <symlink name="fence_ilo3" shortdesc="Fence agent for HP iLO3"/>
  <symlink name="fence_ilo4" shortdesc="Fence agent for HP iLO4"/>
  <symlink name="fence_ilo5" shortdesc="Fence agent for HP iLO5"/>
  <symlink name="fence_imm" shortdesc="Fence agent for IBM Integrated Management Module"/>
  <symlink name="fence_idrac" shortdesc="Fence agent for Dell iDRAC"/>
  <longdesc>fence_ipmilan is an I/O Fencing agent which can be used with machines controlled
  ↳ by IPMI. This agent calls support software ipmitool (http://ipmitool.sf.net/). WARNING! This
  ↳ fence agent might report success before the node is powered off. You should use -m/method
  ↳ onoff if your fence device works correctly with that option.</longdesc>
  <vendor-url/>
  <parameters>
    <parameter name="action" unique="0" required="0">
      <getopt mixed="-o, --action=[action]"/>
      <content type="string" default="reboot"/>
      <shortdesc lang="en">Fencing action</shortdesc>
    </parameter>
    <parameter name="auth" unique="0" required="0">
      <getopt mixed="-A, --auth=[auth]"/>
      <content type="select">
        <option value="md5"/>
        <option value="password"/>
        <option value="none"/>
      </content>
      <shortdesc lang="en">IPMI Lan Auth type.</shortdesc>
    </parameter>
    <parameter name="cipher" unique="0" required="0">
      <getopt mixed="-C, --cipher=[cipher]"/>
      <content type="string"/>
      <shortdesc lang="en">Ciphersuite to use (same as ipmitool -C parameter)</shortdesc>
    </parameter>
    <parameter name="hexadecimal_kg" unique="0" required="0">
      <getopt mixed="--hexadecimal-kg=[key]"/>
      <content type="string"/>
      <shortdesc lang="en">Hexadecimal-encoded Kg key for IPMIv2 authentication</shortdesc>
    </parameter>
    <parameter name="ip" unique="0" required="0" obsoletes="ipaddr">
      <getopt mixed="-a, --ip=[ip]"/>
      <content type="string"/>
      <shortdesc lang="en">IP address or hostname of fencing device</shortdesc>
    </parameter>
    <parameter name="ipaddr" unique="0" required="0" deprecated="1">
      <getopt mixed="-a, --ip=[ip]"/>
      <content type="string"/>
      <shortdesc lang="en">IP address or hostname of fencing device</shortdesc>
    </parameter>
  </parameters>
```

(continues on next page)

(continued from previous page)

```

<parameter name="ipport" unique="0" required="0">
  <getopt mixed="-u, --ipport=[port]"/>
  <content type="integer" default="623"/>
  <shortdesc lang="en">TCP/UDP port to use for connection with device</shortdesc>
</parameter>
<parameter name="lanplus" unique="0" required="0">
  <getopt mixed="-P, --lanplus"/>
  <content type="boolean" default="0"/>
  <shortdesc lang="en">Use Lanplus to improve security of connection</shortdesc>
</parameter>
<parameter name="login" unique="0" required="0" deprecated="1">
  <getopt mixed="-l, --username=[name]"/>
  <content type="string"/>
  <shortdesc lang="en">Login name</shortdesc>
</parameter>
<parameter name="method" unique="0" required="0">
  <getopt mixed="-m, --method=[method]"/>
  <content type="select" default="onoff">
    <option value="onoff"/>
    <option value="cycle"/>
  </content>
  <shortdesc lang="en">Method to fence</shortdesc>
</parameter>
<parameter name="passwd" unique="0" required="0" deprecated="1">
  <getopt mixed="-p, --password=[password]"/>
  <content type="string"/>
  <shortdesc lang="en">Login password or passphrase</shortdesc>
</parameter>
<parameter name="passwd_script" unique="0" required="0" deprecated="1">
  <getopt mixed="-S, --password-script=[script]"/>
  <content type="string"/>
  <shortdesc lang="en">Script to run to retrieve password</shortdesc>
</parameter>
<parameter name="password" unique="0" required="0" obsoletes="passwd">
  <getopt mixed="-p, --password=[password]"/>
  <content type="string"/>
  <shortdesc lang="en">Login password or passphrase</shortdesc>
</parameter>
<parameter name="password_script" unique="0" required="0" obsoletes="passwd_script">
  <getopt mixed="-S, --password-script=[script]"/>
  <content type="string"/>
  <shortdesc lang="en">Script to run to retrieve password</shortdesc>
</parameter>
<parameter name="plug" unique="0" required="0" obsoletes="port">
  <getopt mixed="-n, --plug=[ip]"/>
  <content type="string"/>
  <shortdesc lang="en">IP address or hostname of fencing device (together with --port-as-
↪ip)</shortdesc>
</parameter>
<parameter name="port" unique="0" required="0" deprecated="1">
  <getopt mixed="-n, --plug=[ip]"/>
  <content type="string"/>
  <shortdesc lang="en">IP address or hostname of fencing device (together with --port-as-
↪ip)</shortdesc>
</parameter>
<parameter name="privlvl" unique="0" required="0">

```

(continues on next page)

(continued from previous page)

```

<getopt mixed="-L, --privlvl=[level]"/>
<content type="select" default="administrator">
  <option value="callback"/>
  <option value="user"/>
  <option value="operator"/>
  <option value="administrator"/>
</content>
<shortdesc lang="en">Privilege level on IPMI device</shortdesc>
</parameter>
<parameter name="target" unique="0" required="0">
<getopt mixed="--target=[targetaddress]"/>
<content type="string"/>
<shortdesc lang="en">Bridge IPMI requests to the remote target address</shortdesc>
</parameter>
<parameter name="username" unique="0" required="0" obsoletes="login">
<getopt mixed="-l, --username=[name]"/>
<content type="string"/>
<shortdesc lang="en">Login name</shortdesc>
</parameter>
<parameter name="quiet" unique="0" required="0">
<getopt mixed="-q, --quiet"/>
<content type="boolean"/>
<shortdesc lang="en">Disable logging to stderr. Does not affect --verbose or --debug-
↪file or logging to syslog.</shortdesc>
</parameter>
<parameter name="verbose" unique="0" required="0">
<getopt mixed="-v, --verbose"/>
<content type="boolean"/>
<shortdesc lang="en">Verbose mode</shortdesc>
</parameter>
<parameter name="debug" unique="0" required="0" deprecated="1">
<getopt mixed="-D, --debug-file=[debugfile]"/>
<content type="string"/>
<shortdesc lang="en">Write debug information to given file</shortdesc>
</parameter>
<parameter name="debug_file" unique="0" required="0" obsoletes="debug">
<getopt mixed="-D, --debug-file=[debugfile]"/>
<content type="string"/>
<shortdesc lang="en">Write debug information to given file</shortdesc>
</parameter>
<parameter name="version" unique="0" required="0">
<getopt mixed="-V, --version"/>
<content type="boolean"/>
<shortdesc lang="en">Display version information and exit</shortdesc>
</parameter>
<parameter name="help" unique="0" required="0">
<getopt mixed="-h, --help"/>
<content type="boolean"/>
<shortdesc lang="en">Display help and exit</shortdesc>
</parameter>
<parameter name="delay" unique="0" required="0">
<getopt mixed="--delay=[seconds]"/>
<content type="second" default="0"/>
<shortdesc lang="en">Wait X seconds before fencing is started</shortdesc>
</parameter>
<parameter name="ipmitool_path" unique="0" required="0">

```

(continues on next page)

(continued from previous page)

```

    <getopt mixed="--ipmitool-path=[path]"/>
    <content type="string" default="/usr/bin/ipmitool"/>
    <shortdesc lang="en">Path to ipmitool binary</shortdesc>
  </parameter>
  <parameter name="login_timeout" unique="0" required="0">
    <getopt mixed="--login-timeout=[seconds]"/>
    <content type="second" default="5"/>
    <shortdesc lang="en">Wait X seconds for cmd prompt after login</shortdesc>
  </parameter>
  <parameter name="port_as_ip" unique="0" required="0">
    <getopt mixed="--port-as-ip"/>
    <content type="boolean"/>
    <shortdesc lang="en">Make "port/plug" to be an alias to IP address</shortdesc>
  </parameter>
  <parameter name="power_timeout" unique="0" required="0">
    <getopt mixed="--power-timeout=[seconds]"/>
    <content type="second" default="20"/>
    <shortdesc lang="en">Test X seconds for status change after ON/OFF</shortdesc>
  </parameter>
  <parameter name="power_wait" unique="0" required="0">
    <getopt mixed="--power-wait=[seconds]"/>
    <content type="second" default="2"/>
    <shortdesc lang="en">Wait X seconds after issuing ON/OFF</shortdesc>
  </parameter>
  <parameter name="shell_timeout" unique="0" required="0">
    <getopt mixed="--shell-timeout=[seconds]"/>
    <content type="second" default="3"/>
    <shortdesc lang="en">Wait X seconds for cmd prompt after issuing command</shortdesc>
  </parameter>
  <parameter name="retry_on" unique="0" required="0">
    <getopt mixed="--retry-on=[attempts]"/>
    <content type="integer" default="1"/>
    <shortdesc lang="en">Count of attempts to retry power on</shortdesc>
  </parameter>
  <parameter name="sudo" unique="0" required="0" deprecated="1">
    <getopt mixed="--use-sudo"/>
    <content type="boolean"/>
    <shortdesc lang="en">Use sudo (without password) when calling 3rd party software</
↪shortdesc>
  </parameter>
  <parameter name="use_sudo" unique="0" required="0" obsoletes="sudo">
    <getopt mixed="--use-sudo"/>
    <content type="boolean"/>
    <shortdesc lang="en">Use sudo (without password) when calling 3rd party software</
↪shortdesc>
  </parameter>
  <parameter name="sudo_path" unique="0" required="0">
    <getopt mixed="--sudo-path=[path]"/>
    <content type="string" default="/usr/bin/sudo"/>
    <shortdesc lang="en">Path to sudo binary</shortdesc>
  </parameter>
</parameters>
<actions>
  <action name="on" automatic="0"/>
  <action name="off"/>
  <action name="reboot"/>

```

(continues on next page)

(continued from previous page)

```

<action name="status"/>
<action name="monitor"/>
<action name="metadata"/>
<action name="manpage"/>
<action name="validate-all"/>
<action name="diag"/>
<action name="stop" timeout="20s"/>
<action name="start" timeout="20s"/>
</actions>
</resource-agent>

```

Once we've decided what parameter values we think we need, it is a good idea to run the fence agent's status action manually, to verify that our values work correctly:

```

# fence_ipmilan --lanplus -a 192.0.2.1 -l testuser -p abc123 -o status

Chassis Power is on

```

- Based on that, we might create a fencing resource configuration like this in `stonith.xml` (or any file name, just use the same name with `cibadmin` later):

```

<primitive id="Fencing-pcmk-1" class="stonith" type="fence_ipmilan" >
  <instance_attributes id="Fencing-params" >
    <nvpair id="Fencing-lanplus" name="lanplus" value="1" />
    <nvpair id="Fencing-ip" name="ip" value="192.0.2.1" />
    <nvpair id="Fencing-password" name="password" value="testuser" />
    <nvpair id="Fencing-username" name="username" value="abc123" />
  </instance_attributes>
  <operations >
    <op id="Fencing-monitor-10m" interval="10m" name="monitor" timeout="300s" />
  </operations>
</primitive>

```

Note: Even though the man page shows that the `action` parameter is supported, we do not provide that in the resource configuration. Pacemaker will supply an appropriate action whenever the fence device must be used.

- In this case, we don't need to configure `pcm_k_host_map` because `fence_ipmilan` ignores the target node name and instead uses its `ip` parameter to know how to contact the IPMI controller.
- We do need to let Pacemaker know which cluster node can be fenced by this device, since `fence_ipmilan` doesn't support the `list` action. Add a line like this to the agent's instance attributes:

```

<nvpair id="Fencing-pcmk_host_list" name="pcm_k_host_list" value="pcm_k-1" />

```

- We don't need to configure `pcm_k_host_argument` since `ip` is all the fence agent needs (it ignores the target name).
- Make the configuration active:

```

# cibadmin --create --scope resources --xml-file stonith.xml

```

- Set `stonith-enabled` to `true` (this only has to be done once):

```
# crm_attribute --type crm_config --name stonith-enabled --update true
```

9. Since our cluster is still in testing, we can reboot `pcmk-1` without bothering anyone, so we'll test our fencing configuration by running this from one of the other cluster nodes:

```
# stonith_admin --reboot pcmk-1
```

Then we will verify that the node did, in fact, reboot.

We can repeat that process to create a separate fencing resource for each node.

With some other fence device types, a single fencing resource is able to be used for all nodes. In fact, we could do that with `fence_ipmilan`, using the `port-as-ip` parameter along with `pcmk_host_map`. Either approach is fine.

2.8.15 Fencing Topologies

Pacemaker supports fencing nodes with multiple devices through a feature called *fencing topologies*. Fencing topologies may be used to provide alternative devices in case one fails, or to require multiple devices to all be executed successfully in order to consider the node successfully fenced, or even a combination of the two.

Create the individual devices as you normally would, then define one or more `fencing-level` entries in the `fencing-topology` section of the configuration.

- Each fencing level is attempted in order of ascending `index`. Allowed values are 1 through 9.
- If a device fails, processing terminates for the current level. No further devices in that level are exercised, and the next level is attempted instead.
- If the operation succeeds for all the listed devices in a level, the level is deemed to have passed.
- The operation is finished when a level has passed (success), or all levels have been attempted (failed).
- If the operation failed, the next step is determined by the scheduler and/or the controller.

Some possible uses of topologies include:

- Try on-board IPMI, then an intelligent power switch if that fails
- Try fabric fencing of both disk and network, then fall back to power fencing if either fails
- Wait up to a certain time for a kernel dump to complete, then cut power to the node

Table 20: Attributes of a fencing-level Element

Attribute	Description
<code>id</code>	A unique name for this element (required)
<code>target</code>	The name of a single node to which this level applies
<code>target-pattern</code>	An extended regular expression (as defined in <i>POSIX</i>) matching the names of nodes to which this level applies
<code>target-attribute</code>	The name of a node attribute that is set (to <code>target-value</code>) for nodes to which this level applies
<code>target-value</code>	The node attribute value (of <code>target-attribute</code>) that is set for nodes to which this level applies
<code>index</code>	The order in which to attempt the levels. Levels are attempted in ascending order <i>until one succeeds</i> . Valid values are 1 through 9.
<code>devices</code>	A comma-separated list of devices that must all be tried for this level

Note: Fencing topology with different devices for different nodes

```

<cib crm_feature_set="3.6.0" validate-with="pacemaker-3.5" admin_epoch="1" epoch="0" num_updates="0
↳">
  <configuration>
    ...
    <fencing-topology>
      <!-- For pcmk-1, try poison-pill and fail back to power -->
      <fencing-level id="f-p1.1" target="pcmk-1" index="1" devices="poison-pill"/>
      <fencing-level id="f-p1.2" target="pcmk-1" index="2" devices="power"/>

      <!-- For pcmk-2, try disk and network, and fail back to power -->
      <fencing-level id="f-p2.1" target="pcmk-2" index="1" devices="disk,network"/>
      <fencing-level id="f-p2.2" target="pcmk-2" index="2" devices="power"/>
    </fencing-topology>
    ...
  </configuration>
  <status/>
</cib>

```

Example Dual-Layer, Dual-Device Fencing Topologies

The following example illustrates an advanced use of `fencing-topology` in a cluster with the following properties:

- 2 nodes (prod-mysql1 and prod-mysql2)
- the nodes have IPMI controllers reachable at 192.0.2.1 and 192.0.2.2
- the nodes each have two independent Power Supply Units (PSUs) connected to two independent Power Distribution Units (PDUs) reachable at 198.51.100.1 (port 10 and port 11) and 203.0.113.1 (port 10 and port 11)
- fencing via the IPMI controller uses the `fence_ipmilan` agent (1 fence device per controller, with each device targeting a separate node)
- fencing via the PDUs uses the `fence_apc_snmp` agent (1 fence device per PDU, with both devices targeting both nodes)
- a random delay is used to lessen the chance of a “death match”
- fencing topology is set to try IPMI fencing first then dual PDU fencing if that fails

In a node failure scenario, Pacemaker will first select `fence_ipmilan` to try to kill the faulty node. Using the fencing topology, if that method fails, it will then move on to selecting `fence_apc_snmp` twice (once for the first PDU, then again for the second PDU).

The fence action is considered successful only if both PDUs report the required status. If any of them fails, fencing loops back to the first fencing method, `fence_ipmilan`, and so on, until the node is fenced or the fencing action is cancelled.

Note: First fencing method: single IPMI device per target

Each cluster node has its own dedicated IPMI controller that can be contacted for fencing using the following primitives:

```

<primitive class="stonith" id="fence_prod-mysql1_ipmi" type="fence_ipmilan">
  <instance_attributes id="fence_prod-mysql1_ipmi-instance_attributes">
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-ipaddr" name="ipaddr" value="192.0.2.1"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-passwd" name="passwd" value="finishme"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-lanplus" name="lanplus" value="true"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-pcmk_host_list" name="pcm_k_host_list"
    ↪ value="prod-mysql1"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max"
    ↪ value="8s"/>
  </instance_attributes>
</primitive>
<primitive class="stonith" id="fence_prod-mysql2_ipmi" type="fence_ipmilan">
  <instance_attributes id="fence_prod-mysql2_ipmi-instance_attributes">
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-ipaddr" name="ipaddr" value="192.0.2.2"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-passwd" name="passwd" value="finishme"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-lanplus" name="lanplus" value="true"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-pcmk_host_list" name="pcm_k_host_list"
    ↪ value="prod-mysql2"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max"
    ↪ value="8s"/>
  </instance_attributes>
</primitive>

```

Note: Second fencing method: dual PDU devices

Each cluster node also has 2 distinct power supplies controlled by 2 distinct PDUs:

- Node 1: PDU 1 port 10 and PDU 2 port 10
- Node 2: PDU 1 port 11 and PDU 2 port 11

The matching fencing agents are configured as follows:

```

<primitive class="stonith" id="fence_apc1" type="fence_apc_snmp">
  <instance_attributes id="fence_apc1-instance_attributes">
    <nvpair id="fence_apc1-instance_attributes-ipaddr" name="ipaddr" value="198.51.100.1"/>
    <nvpair id="fence_apc1-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_apc1-instance_attributes-passwd" name="passwd" value="fencing"/>
    <nvpair id="fence_apc1-instance_attributes-pcmk_host_list"
    ↪ name="pcm_k_host_map" value="prod-mysql1:10;prod-mysql2:11"/>
    <nvpair id="fence_apc1-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max" value="8s"/>
  </instance_attributes>
</primitive>
<primitive class="stonith" id="fence_apc2" type="fence_apc_snmp">
  <instance_attributes id="fence_apc2-instance_attributes">
    <nvpair id="fence_apc2-instance_attributes-ipaddr" name="ipaddr" value="203.0.113.1"/>
    <nvpair id="fence_apc2-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_apc2-instance_attributes-passwd" name="passwd" value="fencing"/>
    <nvpair id="fence_apc2-instance_attributes-pcmk_host_list"
    ↪ name="pcm_k_host_map" value="prod-mysql1:10;prod-mysql2:11"/>
    <nvpair id="fence_apc2-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max" value="8s"/>
  </instance_attributes>
</primitive>

```

Note: Fencing topology

Now that all the fencing resources are defined, it's time to create the right topology. We want to first fence using IPMI and if that does not work, fence both PDUs to effectively and surely kill the node.

```
<fencing-topology>
  <fencing-level id="level-1-1" target="prod-mysql1" index="1" devices="fence_prod-mysql1_ipmi" />
  <fencing-level id="level-1-2" target="prod-mysql1" index="2" devices="fence_apc1,fence_apc2" />
  <fencing-level id="level-2-1" target="prod-mysql2" index="1" devices="fence_prod-mysql2_ipmi" />
  <fencing-level id="level-2-2" target="prod-mysql2" index="2" devices="fence_apc1,fence_apc2" />
</fencing-topology>
```

In `fencing-topology`, the lowest `index` value for a target determines its first fencing method.

2.8.16 Remapping Reboots

When the cluster needs to reboot a node, whether because `stonith-action` is `reboot` or because a reboot was requested externally (such as by `stonith_admin --reboot`), it will remap that to other commands in two cases:

- If the chosen fencing device does not support the `reboot` command, the cluster will ask it to perform `off` instead.
- If a fencing topology level with multiple devices must be executed, the cluster will ask all the devices to perform `off`, then ask the devices to perform `on`.

To understand the second case, consider the example of a node with redundant power supplies connected to intelligent power switches. Rebooting one switch and then the other would have no effect on the node. Turning both switches off, and then on, actually reboots the node.

In such a case, the fencing operation will be treated as successful as long as the `off` commands succeed, because then it is safe for the cluster to recover any resources that were on the node. Timeouts and errors in the `on` phase will be logged but ignored.

When a reboot operation is remapped, any action-specific timeout for the remapped action will be used (for example, `pcmk_off_timeout` will be used when executing the `off` command, not `pcmk_reboot_timeout`).

2.9 Collective Resources

Pacemaker supports several types of *collective* resources, which consist of multiple, related resource instances.

2.9.1 Groups - A Syntactic Shortcut

One of the most common elements of a cluster is a set of resources that need to be located together, start sequentially, and stop in the reverse order. To simplify this configuration, we support the concept of groups.

A group of two primitive resources

```

<group id="shortcut">
  <primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
    <instance_attributes id="params-public-ip">
      <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
    </instance_attributes>
  </primitive>
  <primitive id="Email" class="systemd" type="exim"/>
</group>

```

Although the example above contains only two resources, there is no limit to the number of resources a group can contain. The example is also sufficient to explain the fundamental properties of a group:

- Resources are started in the order they appear in (**Public-IP** first, then **Email**)
- Resources are stopped in the reverse order to which they appear in (**Email** first, then **Public-IP**)

If a resource in the group can't run anywhere, then nothing after that is allowed to run, too.

- If **Public-IP** can't run anywhere, neither can **Email**;
- but if **Email** can't run anywhere, this does not affect **Public-IP** in any way

The group above is logically equivalent to writing:

How the cluster sees a group resource

```

<configuration>
  <resources>
    <primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
      <instance_attributes id="params-public-ip">
        <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
      </instance_attributes>
    </primitive>
    <primitive id="Email" class="systemd" type="exim"/>
  </resources>
  <constraints>
    <rsc_colocation id="xxx" rsc="Email" with-rsc="Public-IP" score="INFINITY"/>
    <rsc_order id="yyy" first="Public-IP" then="Email"/>
  </constraints>
</configuration>

```

Obviously as the group grows bigger, the reduced configuration effort can become significant.

Another (typical) example of a group is a DRBD volume, the filesystem mount, an IP address, and an application that uses them.

Group Properties

Table 21: Properties of a Group Resource

Field	Description
id	A unique name for the group
description	Arbitrary text for user's use (ignored by Pacemaker)

Group Options

Groups inherit the `priority`, `target-role`, and `is-managed` properties from primitive resources. See *Resource Options* for information about those properties.

Group Instance Attributes

Groups have no instance attributes. However, any that are set for the group object will be inherited by the group's children.

Group Contents

Groups may only contain a collection of cluster resources (see *Resource Properties*). To refer to a child of a group resource, just use the child's `id` instead of the group's.

Group Constraints

Although it is possible to reference a group's children in constraints, it is usually preferable to reference the group itself.

Some constraints involving groups

```
<constraints>
  <rsc_location id="group-prefers-node1" rsc="shortcut" node="node1" score="500"/>
  <rsc_colocation id="webserver-with-group" rsc="Webserver" with-rsc="shortcut"/>
  <rsc_order id="start-group-then-webserver" first="Webserver" then="shortcut"/>
</constraints>
```

Group Stickiness

Stickiness, the measure of how much a resource wants to stay where it is, is additive in groups. Every active resource of the group will contribute its stickiness value to the group's total. So if the default `resource-stickiness` is 100, and a group has seven members, five of which are active, then the group as a whole will prefer its current location with a score of 500.

2.9.2 Clones - Resources That Can Have Multiple Active Instances

Clone resources are resources that can have more than one copy active at the same time. This allows you, for example, to run a copy of a daemon on every node. You can clone any primitive or group resource¹.

Anonymous versus Unique Clones

A clone resource is configured to be either *anonymous* or *globally unique*.

Anonymous clones are the simplest. These behave completely identically everywhere they are running. Because of this, there can be only one instance of an anonymous clone active per node.

¹ Of course, the service must support running multiple instances.

The instances of globally unique clones are distinct entities. All instances are launched identically, but one instance of the clone is not identical to any other instance, whether running on the same node or a different node. As an example, a cloned IP address can use special kernel functionality such that each instance handles a subset of requests for the same IP address.

Promotable clones

If a clone is *promotable*, its instances can perform a special role that Pacemaker will manage via the `promote` and `demote` actions of the resource agent.

Services that support such a special role have various terms for the special role and the default role: primary and secondary, master and replica, controller and worker, etc. Pacemaker uses the terms *promoted* and *unpromoted* to be agnostic to what the service calls them or what they do.

All that Pacemaker cares about is that an instance comes up in the unpromoted role when started, and the resource agent supports the `promote` and `demote` actions to manage entering and exiting the promoted role.

Clone Properties

Table 22: Properties of a Clone Resource

Field	Description
id	A unique name for the clone
description	Arbitrary text for user's use (ignored by Pacemaker)

Clone Options

Options inherited from primitive resources: `priority`, `target-role`, `is-managed`

Table 23: Clone-specific configuration options

Field	Default	Description
globally-unique	true if <code>clone-node-max</code> is greater than 1 (<i>since 3.0.0</i>), otherwise false	If true , each clone instance performs a distinct function, such that a single node can run more than one instance at the same time
clone-max	0	The maximum number of clone instances that can be started across the entire cluster. If 0, the number of nodes in the cluster will be used.
clone-node-max	1	If the clone is globally unique, this is the maximum number of clone instances that can be started on a single node
clone-min	0	Require at least this number of clone instances to be runnable before allowing resources depending on the clone to be runnable. A value of 0 means require all clone instances to be runnable.
notify	false	Call the resource agent's notify action for all active instances, before and after starting or stopping any clone instance. The resource agent must support this action. Allowed values: false , true

Continued on next page

Table 23 – continued from previous page

Field	Default	Description
ordered	false	If true , clone instances must be started sequentially instead of in parallel. Allowed values: false , true
interleave	false	When this clone is ordered relative to another clone, if this option is false (the default), the ordering is relative to <i>all</i> instances of the other clone, whereas if this option is true , the ordering is relative only to instances on the same node. Allowed values: false , true
promotable	false	If true , clone instances can perform a special role that Pacemaker will manage via the resource agent's promote and demote actions. The resource agent must support these actions. Allowed values: false , true
promoted-max	1	If promotable is true , the number of instances that can be promoted at one time across the entire cluster
promoted-node-max	1	If the clone is promotable and globally unique, this is the number of instances that can be promoted at one time on a single node (up to <code>clone-node-max</code>)

Note: Deprecated Terminology

In older documentation and online examples, you may see promotable clones referred to as *multi-state*, *stateful*, or *master/slave*; these mean the same thing as *promotable*. Certain syntax is supported for backward compatibility, but is deprecated and will be removed in a future version:

- Using the `master-max` meta-attribute instead of `promoted-max`
- Using the `master-node-max` meta-attribute instead of `promoted-node-max`
- Using `Master` as a role name instead of `Promoted`
- Using `Slave` as a role name instead of `Unpromoted`

Clone Contents

Clones must contain exactly one primitive or group resource.

A clone that runs a web server on all nodes

```
<clone id="apache-clone">
  <primitive id="apache" class="systemd" type="httpd">
    <operations>
      <op id="apache-monitor" name="monitor" interval="30"/>
    </operations>
  </primitive>
</clone>
```

Warning: You should never reference the name of a clone's child (the primitive or group resource being cloned). If you think you need to do this, you probably need to re-evaluate your design.

Clone Instance Attribute

Clones have no instance attributes; however, any that are set here will be inherited by the clone's child.

Clone Constraints

In most cases, a clone will have a single instance on each active cluster node. If this is not the case, you can indicate which nodes the cluster should preferentially assign copies to with resource location constraints. These constraints are written no differently from those for primitive resources except that the clone's `id` is used.

Some constraints involving clones

```
<constraints>
  <rsc_location id="clone-prefers-node1" rsc="apache-clone" node="node1" score="500"/>
  <rsc_colocation id="stats-with-clone" rsc="apache-stats" with="apache-clone"/>
  <rsc_order id="start-clone-then-stats" first="apache-clone" then="apache-stats"/>
</constraints>
```

Ordering constraints behave slightly differently for clones. In the example above, `apache-stats` will wait until all copies of `apache-clone` that need to be started have done so before being started itself. Only if *no* copies can be started will `apache-stats` be prevented from being active. Additionally, the clone will wait for `apache-stats` to be stopped before stopping itself.

Colocation of a primitive or group resource with a clone means that the resource can run on any node with an active instance of the clone. The cluster will choose an instance based on where the clone is running and the resource's own location preferences.

Colocation between clones is also possible. If one clone **A** is colocated with another clone **B**, the set of allowed locations for **A** is limited to nodes on which **B** is (or will be) active. Placement is then performed normally.

Promotable Clone Constraints

For promotable clone resources, the `first-action` and/or `then-action` fields for ordering constraints may be set to `promote` or `demote` to constrain the promoted role, and colocation constraints may contain `rsc-role` and/or `with-rsc-role` fields.

Constraints involving promotable clone resources

```
<constraints>
  <rsc_location id="db-prefers-node1" rsc="database" node="node1" score="500"/>
  <rsc_colocation id="backup-with-db-unpromoted" rsc="backup"
    with-rsc="database" with-rsc-role="Unpromoted"/>
  <rsc_colocation id="myapp-with-db-promoted" rsc="myApp"
    with-rsc="database" with-rsc-role="Promoted"/>
  <rsc_order id="start-db-before-backup" first="database" then="backup"/>
  <rsc_order id="promote-db-then-app" first="database" first-action="promote"
    then="myApp" then-action="start"/>
</constraints>
```

In the example above, **myApp** will wait until one of the database copies has been started and promoted before being started itself on the same node. Only if no copies can be promoted will **myApp** be prevented from being active. Additionally, the cluster will wait for **myApp** to be stopped before demoting the database.

Colocation of a primitive or group resource with a promotable clone resource means that it can run on any node with an active instance of the promotable clone resource that has the specified role (**Promoted** or **Unpromoted**). In the example above, the cluster will choose a location based on where database is running in the promoted role, and if there are multiple promoted instances it will also factor in **myApp**'s own location preferences when deciding which location to choose.

Colocation with regular clones and other promotable clone resources is also possible. In such cases, the set of allowed locations for the **rsc** clone is (after role filtering) limited to nodes on which the **with-rsc** promotable clone resource is (or will be) in the specified role. Placement is then performed as normal.

Using Promotable Clone Resources in Colocation Sets

When a promotable clone is used in a *resource set* inside a colocation constraint, the resource set may take a **role** attribute.

In the following example, an instance of **B** may be promoted only on a node where **A** is in the promoted role. Additionally, resources **C** and **D** must be located on a node where both **A** and **B** are promoted.

Colocate C and D with A's and B's promoted instances

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-example-1" sequential="true" role="Promoted">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="colocated-set-example-2" sequential="true">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

Using Promotable Clone Resources in Ordered Sets

When a promotable clone is used in a *resource set* inside an ordering constraint, the resource set may take an **action** attribute.

Start C and D after first promoting A and B

```
<constraints>
  <rsc_order id="order-1" score="INFINITY" >
    <resource_set id="ordered-set-1" sequential="true" action="promote">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="true" action="start">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_order>
</constraints>
```

In the above example, **B** cannot be promoted until **A** has been promoted. Additionally, resources **C** and **D** must wait until **A** and **B** have been promoted before they can start.

Clone Stickiness

To achieve stable assignments, clones are slightly sticky by default. If no value for `resource-stickiness` is provided, the clone will use a value of 1. Being a small value, it causes minimal disturbance to the score calculations of other resources but is enough to prevent Pacemaker from needlessly moving instances around the cluster.

Note: For globally unique clones, this may result in multiple instances of the clone staying on a single node, even after another eligible node becomes active (for example, after being put into standby mode then made active again). If you do not want this behavior, specify a `resource-stickiness` of 0 for the clone temporarily and let the cluster adjust, then set it back to 1 if you want the default behavior to apply again.

Important: If `resource-stickiness` is set in the `rsc_defaults` section, it will apply to clone instances as well. This means an explicit `resource-stickiness` of 0 in `rsc_defaults` works differently from the implicit default used when `resource-stickiness` is not specified.

Monitoring Promotable Clone Resources

The usual monitor actions are insufficient to monitor a promotable clone resource, because Pacemaker needs to verify not only that the resource is active, but also that its actual role matches its intended one.

Define two monitoring actions: the usual one will cover the unpromoted role, and an additional one with `role="Promoted"` will cover the promoted role.

Monitoring both states of a promotable clone resource

```
<clone id="myPromotableRsc">
  <meta_attributes id="myPromotableRsc-meta">
    <nvpair name="promotable" value="true"/>
  </meta_attributes>
  <primitive id="myRsc" class="ocf" type="myApp" provider="myCorp">
    <operations>
      <op id="public-ip-unpromoted-check" name="monitor" interval="60"/>
      <op id="public-ip-promoted-check" name="monitor" interval="61" role="Promoted"/>
    </operations>
  </primitive>
</clone>
```

Important: It is crucial that *every* monitor operation has a different interval! Pacemaker currently differentiates between operations only by resource and interval; so if (for example) a promotable clone resource had the same monitor interval for both roles, Pacemaker would ignore the role when checking the status – which would cause unexpected return codes, and therefore unnecessary complications.

Determining Which Instance is Promoted

Pacemaker can choose a promotable clone instance to be promoted in one of two ways:

- Promotion scores: These are node attributes set via the `crm_attribute` command using the `--promotion` option, which generally would be called by the resource agent's start action if it supports promotable clones. This tool automatically detects both the resource and host, and should be used to set a preference for being promoted. Based on this, `promoted-max`, and `promoted-node-max`, the instance(s) with the highest preference will be promoted.
- Constraints: Location constraints can indicate which nodes are most preferred to be promoted.

Explicitly preferring node1 to be promoted

```
<rsc_location id="promoted-location" rsc="myPromotableRsc">
  <rule id="promoted-rule" score="100" role="Promoted">
    <expression id="promoted-exp" attribute="#uname" operation="eq" value="node1"/>
  </rule>
</rsc_location>
```

2.9.3 Bundles - Containerized Resources

Pacemaker supports a special syntax for launching a service inside a `container` with any infrastructure it requires: the *bundle*.

Pacemaker bundles support `Docker` and `podman` (since 2.0.1) container technologies.²

A bundle for a containerized web server

```
<bundle id="httpd-bundle">
  <podman image="pcmk:http" replicas="3"/>
  <network ip-range-start="192.168.122.131"
    host-netmask="24"
    host-interface="eth0">
    <port-mapping id="httpd-port" port="80"/>
  </network>
  <storage>
    <storage-mapping id="httpd-syslog"
      source-dir="/dev/log"
      target-dir="/dev/log"
      options="rw"/>
    <storage-mapping id="httpd-root"
      source-dir="/srv/html"
      target-dir="/var/www/html"
      options="rw,Z"/>
    <storage-mapping id="httpd-logs"
      source-dir-root="/var/log/pacemaker/bundles"
      target-dir="/etc/httpd/logs"
      options="rw,Z"/>
  </storage>
  <primitive class="ocf" id="httpd" provider="heartbeat" type="apache"/>
</bundle>
```

² Docker is a trademark of Docker, Inc. No endorsement by or association with Docker, Inc. is implied.

Bundle Prerequisites

Before configuring a bundle in Pacemaker, the user must install the appropriate container launch technology (Docker or podman), and supply a fully configured container image, on every node allowed to run the bundle.

Pacemaker will create an implicit resource of type **ocf:heartbeat:docker** or **ocf:heartbeat:podman** to manage a bundle's container. The user must ensure that the appropriate resource agent is installed on every node allowed to run the bundle.

Bundle Properties

Table 24: XML Attributes of a bundle Element

Field	Description
id	A unique name for the bundle (required)
description	Arbitrary text for user's use (ignored by Pacemaker)

A bundle must contain exactly one `docker` or `podman` element.

Bundle Container Properties

Table 25: XML attributes of a docker or podman Element

Attribute	Default	Description
image		Container image tag (required)
replicas	Value of <code>promoted-max</code> if that is positive, else 1	A positive integer specifying the number of container instances to launch
replicas-per-host	1	A positive integer specifying the number of container instances allowed to run on a single node
promoted-max	0	A non-negative integer that, if positive, indicates that the containerized service should be treated as a promotable service, with this many replicas allowed to run the service in the promoted role
network		If specified, this will be passed to the <code>docker run</code> or <code>podman run</code> command as the network setting for the container.
run-command	<code>/usr/sbin/pacemaker-remoted</code> if bundle contains a primitive , otherwise none	This command will be run inside the container when launching it ("PID 1"). If the bundle contains a primitive , this command <i>must</i> start <code>pacemaker-remoted</code> (but could, for example, be a script that does other stuff, too).
options		Extra command-line options to pass to the <code>docker run</code> or <code>podman run</code> command

Note: Considerations when using cluster configurations or container images from Pacemaker 1.1:

- If the container image has a pre-2.0.0 version of Pacemaker, set `run-command` to `/usr/sbin/`

`pacemaker_remoted` (note the underbar instead of dash).

- `masters` is accepted as an alias for `promoted-max`, but is deprecated since 2.0.0, and support for it will be removed in a future version.

Bundle Network Properties

A bundle may optionally contain one `<network>` element.

Table 26: XML attributes of a `network` Element

Attribute	Default	Description
<code>add-host</code>	TRUE	If TRUE, and <code>ip-range-start</code> is used, Pacemaker will automatically ensure that <code>/etc/hosts</code> inside the containers has entries for each <i>replica name</i> and its assigned IP.
<code>ip-range-start</code>		If specified, Pacemaker will create an implicit <code>ocf:heartbeat:IPaddr2</code> resource for each container instance, starting with this IP address, using up to <code>replicas</code> sequential addresses. These addresses can be used from the host's network to reach the service inside the container, though it is not visible within the container itself. Only IPv4 addresses are currently supported.
<code>host-netmask</code>	32	If <code>ip-range-start</code> is specified, the IP addresses are created with this CIDR netmask (as a number of bits).
<code>host-interface</code>		If <code>ip-range-start</code> is specified, the IP addresses are created on this host interface (by default, it will be determined from the IP address).
<code>control-port</code>	3121	If the bundle contains a <code>primitive</code> , the cluster will use this integer TCP port for communication with Pacemaker Remote inside the container. Changing this is useful when the container is unable to listen on the default port, for example, when the container uses the host's network rather than <code>ip-range-start</code> (in which case <code>replicas-per-host</code> must be 1), or when the bundle may run on a Pacemaker Remote node that is already listening on the default port. Any <code>PCMK_remote_port</code> environment variable set on the host or in the container is ignored for bundle connections.

Note: Replicas are named by the bundle id plus a dash and an integer counter starting with zero. For example, if a bundle named `httpd-bundle` has `replicas=2`, its containers will be named `httpd-bundle-0` and `httpd-bundle-1`.

Additionally, a `network` element may optionally contain one or more `port-mapping` elements.

Table 27: Attributes of a port-mapping Element

Attribute	Default	Description
id		A unique name for the port mapping (required)
port		If this is specified, connections to this TCP port number on the host network (on the container's assigned IP address, if <code>ip-range-start</code> is specified) will be forwarded to the container network. Exactly one of <code>port</code> or <code>range</code> must be specified in a <code>port-mapping</code> .
internal-port	value of <code>port</code>	If <code>port</code> and this are specified, connections to <code>port</code> on the host's network will be forwarded to this port on the container network.
range		If this is specified, connections to these TCP port numbers (expressed as <code>first_port-last_port</code>) on the host network (on the container's assigned IP address, if <code>ip-range-start</code> is specified) will be forwarded to the same ports in the container network. Exactly one of <code>port</code> or <code>range</code> must be specified in a <code>port-mapping</code> .

Note: If the bundle contains a `primitive`, Pacemaker will automatically map the `control-port`, so it is not necessary to specify that port in a `port-mapping`.

Bundle Storage Properties

A bundle may optionally contain one `storage` element. A `storage` element has no properties of its own, but may contain one or more `storage-mapping` elements.

Table 28: Attributes of a storage-mapping Element

Attribute	Default	Description
id		A unique name for the storage mapping (required)
source-dir		The absolute path on the host's filesystem that will be mapped into the container. Exactly one of <code>source-dir</code> and <code>source-dir-root</code> must be specified in a <code>storage-mapping</code> .
source-dir-root		The start of a path on the host's filesystem that will be mapped into the container, using a different subdirectory on the host for each container instance. The subdirectory will be named the same as the <i>replica name</i> . Exactly one of <code>source-dir</code> and <code>source-dir-root</code> must be specified in a <code>storage-mapping</code> .
target-dir		The path name within the container where the host storage will be mapped (required)
options		A comma-separated list of file system mount options to use when mapping the storage

Note: Pacemaker does not define the behavior if the source directory does not already exist on the host. However, it is expected that the container technology and/or its resource agent will create the source

directory in that case.

Note: If the bundle contains a `primitive`, Pacemaker will automatically map the equivalent of `source-dir=/etc/pacemaker/authkey target-dir=/etc/pacemaker/authkey` and `source-dir-root=/var/log/pacemaker/bundles target-dir=/var/log` into the container, so it is not necessary to specify those paths in a `storage-mapping`.

Important: The `PCMK_authkey_location` environment variable must not be set to anything other than the default of `/etc/pacemaker/authkey` on any node in the cluster.

Important: If SELinux is used in enforcing mode on the host, you must ensure the container is allowed to use any storage you mount into it. For Docker and podman bundles, adding “Z” to the mount options will create a container-specific label for the mount that allows the container access.

Bundle Primitive

A bundle may optionally contain one *primitive* resource. The primitive may have operations, instance attributes, and meta-attributes defined, as usual.

If a bundle contains a primitive resource, the container image must include the Pacemaker Remote daemon, and at least one of `ip-range-start` or `control-port` must be configured in the bundle. Pacemaker will create an implicit `ocf:pacemaker:remote` resource for the connection, launch Pacemaker Remote within the container, and monitor and manage the primitive resource via Pacemaker Remote.

If the bundle has more than one container instance (replica), the primitive resource will function as an implicit *clone* – a *promotable clone* if the bundle has `promoted-max` greater than zero.

Note: If you want to pass environment variables to a bundle’s Pacemaker Remote connection or primitive, you have two options:

- Environment variables whose value is the same regardless of the underlying host may be set using the container element’s `options` attribute.
 - If you want variables to have host-specific values, you can use the *storage-mapping* element to map a file on the host as `/etc/pacemaker/pcmk-init.env` in the container (*since 2.0.3*). Pacemaker Remote will parse this file as a shell-like format, with variables set as `NAME=VALUE`, ignoring blank lines and comments starting with “#”.
-

Important: When a bundle has a `primitive`, Pacemaker on all cluster nodes must be able to contact Pacemaker Remote inside the bundle’s containers.

- The containers must have an accessible network (for example, `network` should not be set to “none” with a `primitive`).
- The default, using a distinct network space inside the container, works in combination with `ip-range-start`. Any firewall must allow access from all cluster nodes to the `control-port` on the container IPs.

- If the container shares the host's network space (for example, by setting `network` to "host"), a unique `control-port` should be specified for each bundle. Any firewall must allow access from all cluster nodes to the `control-port` on all cluster and remote node IPs.
-

Bundle Node Attributes

If the bundle has a `primitive`, the primitive's resource agent may want to set node attributes such as *promotion scores*. However, with containers, it is not apparent which node should get the attribute.

If the container uses shared storage that is the same no matter which node the container is hosted on, then it is appropriate to use the promotion score on the bundle node itself.

On the other hand, if the container uses storage exported from the underlying host, then it may be more appropriate to use the promotion score on the underlying host.

Since this depends on the particular situation, the `container-attribute-target` resource meta-attribute allows the user to specify which approach to use. If it is set to `host`, then user-defined node attributes will be checked on the underlying host. If it is anything else, the local node (in this case the bundle node) is used as usual.

This only applies to user-defined attributes; the cluster will always check the local node for cluster-defined attributes such as `#uname`.

If `container-attribute-target` is `host`, the cluster will pass additional environment variables to the primitive's resource agent that allow it to set node attributes appropriately: `CRM_meta_container_attribute_target` (identical to the meta-attribute value) and `CRM_meta_physical_host` (the name of the underlying host).

Note: When called by a resource agent, the `attrd_updater` and `crm_attribute` commands will automatically check those environment variables and set attributes appropriately.

Bundle Meta-Attributes

Any meta-attribute set on a bundle will be inherited by the bundle's primitive and any resources implicitly created by Pacemaker for the bundle.

This includes options such as `priority`, `target-role`, and `is-managed`. See *Resource Options* for more information.

Bundles support clone meta-attributes including `notify`, `ordered`, and `interleave`.

Limitations of Bundles

Restarting pacemaker while a bundle is unmanaged or the cluster is in maintenance mode may cause the bundle to fail.

Bundles may not be explicitly cloned or included in groups. This includes the bundle's primitive and any resources implicitly created by Pacemaker for the bundle. (If `replicas` is greater than 1, the bundle will behave like a clone implicitly.)

Bundles do not have instance attributes, utilization attributes, or operations, though a bundle's primitive may have them.

A bundle with a primitive can run on a Pacemaker Remote node only if the bundle uses a distinct `control-port`.

2.10 Utilization and Placement Strategy

Pacemaker decides where a resource should run by assigning a score to every node, considering factors such as the resource's constraints and stickiness, then assigning the resource to the node with the highest score.

If more than one node has the highest score, Pacemaker by default chooses the one with the least number of assigned resources, or if that is also the same, the one listed first in the CIB. This results in simple load balancing.

Sometimes, simple load balancing is insufficient. Different resources can use significantly different amounts of a node's memory, CPU, and other capacities. Some combinations of resources may strain a node's capacity, causing them to fail or have degraded performance. Or, an administrator may prefer to concentrate resources rather than balance them, to minimize energy consumption by spare nodes.

Pacemaker offers flexibility by allowing you to configure *utilization attributes* specifying capacities that each node provides and each resource requires, as well as a *placement strategy*.

2.10.1 Utilization attributes

You can define any number of utilization attributes to represent capacities of interest (CPU, memory, I/O bandwidth, etc.). Their values must be integers.

The nature and units of the capacities are irrelevant to Pacemaker. It just makes sure that each node has sufficient capacity to run the resources assigned to it.

Specifying CPU and RAM capacities of two nodes

```
<node id="node1" type="normal" uname="node1">
  <utilization id="node1-utilization">
    <nvpair id="node1-utilization-cpu" name="cpu" value="2"/>
    <nvpair id="node1-utilization-memory" name="memory" value="2048"/>
  </utilization>
</node>
<node id="node2" type="normal" uname="node2">
  <utilization id="node2-utilization">
    <nvpair id="node2-utilization-cpu" name="cpu" value="4"/>
    <nvpair id="node2-utilization-memory" name="memory" value="4096"/>
  </utilization>
</node>
```

Specifying CPU and RAM consumed by several resources

```

<primitive id="rsc-small" class="ocf" provider="pacemaker" type="Dummy">
  <utilization id="rsc-small-utilization">
    <nvpair id="rsc-small-utilization-cpu" name="cpu" value="1"/>
    <nvpair id="rsc-small-utilization-memory" name="memory" value="1024"/>
  </utilization>
</primitive>
<primitive id="rsc-medium" class="ocf" provider="pacemaker" type="Dummy">
  <utilization id="rsc-medium-utilization">
    <nvpair id="rsc-medium-utilization-cpu" name="cpu" value="2"/>
    <nvpair id="rsc-medium-utilization-memory" name="memory" value="2048"/>
  </utilization>
</primitive>
<primitive id="rsc-large" class="ocf" provider="pacemaker" type="Dummy">
  <utilization id="rsc-large-utilization">
    <nvpair id="rsc-large-utilization-cpu" name="cpu" value="3"/>
    <nvpair id="rsc-large-utilization-memory" name="memory" value="3072"/>
  </utilization>
</primitive>

```

Utilization attributes for a node may be permanent or (*since 2.1.6*) transient. Permanent attributes persist after Pacemaker is restarted, while transient attributes do not.

Transient utilization attribute for node cluster-1

```

<transient_attributes id="cluster-1">
  <utilization id="status-cluster-1">
    <nvpair id="status-cluster-1-cpu" name="cpu" value="1"/>
  </utilization>
</transient_attributes>

```

Utilization attributes may be configured only on primitive resources. Pacemaker will consider a collective resource's utilization based on the primitives it contains.

Note: Utilization is supported for bundles (*since 2.1.3*), but only for bundles with an inner primitive.

2.10.2 Placement Strategy

The `placement-strategy` cluster option determines how utilization attributes are used. Its allowed values are:

- **default:** The cluster ignores utilization values, and places resources according to (from highest to lowest precedence) assignment scores, the number of resources already assigned to each node, and the order nodes are listed in the CIB.
- **utilization:** The cluster uses the same method as the default strategy to assign a resource to a node, but only nodes with sufficient free capacity to meet the resource's requirements are eligible.
- **balanced:** Only nodes with sufficient free capacity are eligible to run a resource, and the cluster load-balances based on the sum of resource utilization values rather than the number of resources.
- **minimal:** Only nodes with sufficient free capacity are eligible to run a resource, and the cluster concentrates resources on as few nodes as possible.

To look at it another way, when deciding where to run a resource, the cluster starts by considering all nodes, then applies these criteria one by one until a single node remains:

- If `placement-strategy` is `utilization`, `balanced`, or `minimal`, consider only nodes that have sufficient spare capacities to meet the resource's requirements.
- Consider only nodes with the highest score for the resource. Scores take into account factors such as the node's health; the resource's stickiness, failure count on the node, and migration threshold; and constraints.
- If `placement-strategy` is `balanced`, consider only nodes with the most free capacity.
- If `placement-strategy` is `default`, `utilization`, or `balanced`, consider only nodes with the least number of assigned resources.
- If more than one node is eligible after considering all other criteria, choose the one listed first in the CIB.

2.10.3 How Multiple Capacities Combine

If only one type of utilization attribute has been defined, free capacity is a simple numeric comparison.

If multiple utilization attributes have been defined, then the node that has the highest value in the most attribute types has the most free capacity.

For example:

- If `nodeA` has more free `cpus`, and `nodeB` has more free `memory`, then their free capacities are equal.
- If `nodeA` has more free `cpus`, while `nodeB` has more free `memory` and `storage`, then `nodeB` has more free capacity.

2.10.4 Order of Resource Assignment

When assigning resources to nodes, the cluster chooses the next one to assign by considering the following criteria one by one until a single resource is selected:

- Assign the resource with the highest *priority*.
- If any resources are already active, assign the one with the highest score on its current node. This avoids unnecessary resource shuffling.
- Assign the resource with the highest score on its preferred node.
- If more than one resource remains after considering all other criteria, assign the one of them that is listed first in the CIB.

Note: For bundles, only the priority set for the bundle itself matters. If the bundle contains a primitive, the primitive's priority is ignored.

2.10.5 Limitations

The type of problem Pacemaker is dealing with here is known as the [knapsack problem](#) and falls into the [NP-complete](#) category of computer science problems – a fancy way of saying “it takes a really long time to solve”.

In a high-availability cluster, it is unacceptable to spend minutes, let alone hours or days, finding an optimal solution while services are down.

Instead of trying to solve the problem completely, Pacemaker uses a “best effort” algorithm. This arrives at a quick solution, but at the cost of possibly leaving some resources stopped unnecessarily.

Using the example configuration at the start of this chapter, and the balanced placement strategy:

- `rsc-small` would be assigned to `node1`
- `rsc-medium` would be assigned to `node2`
- `rsc-large` would remain inactive

That is not ideal. There are various approaches to dealing with the limitations of Pacemaker’s placement strategy:

- **Ensure you have sufficient physical capacity.**

It might sound obvious, but if the physical capacity of your nodes is maxed out even under normal conditions, failover isn’t going to go well. Even without the utilization feature, you’ll start hitting timeouts and getting secondary failures.

- **Build some buffer into the capacities advertised by the nodes.**

Advertise slightly more resources than we physically have, on the (usually valid) assumption that resources will not always use 100% of their configured utilization. This practice is sometimes called *overcommitting*.

- **Specify resource priorities.**

If the cluster is going to sacrifice services, it should be the ones you care about the least.

2.11 Rules

Rules make a configuration more dynamic, allowing values to depend on conditions such as time of day or the value of a node attribute. For example, rules can:

- Set a higher value for *resource-stickiness* during working hours to minimize downtime, and a lower value on weekends to allow resources to move to their most preferred locations when people aren’t around
- Automatically place the cluster into maintenance mode during a scheduled maintenance window
- Restrict a particular department’s resources to run on certain nodes, as determined by custom resource meta-attributes and node attributes

2.11.1 Rule Options

Each context that supports rules may contain a single `rule` element.

Table 29: Attributes of a rule Element

Name	Type	Default	Description
<code>id</code>	<i>id</i>		A unique name for this element (required)

Continued on next page

Table 29 – continued from previous page

Name	Type	Default	Description
boolean-op	<i>enumeration</i>	and	How to combine conditions if this rule contains more than one. Allowed values: <ul style="list-style-type: none"> • and: the rule is satisfied only if all conditions are satisfied • or: the rule is satisfied if any condition is satisfied

2.11.2 Rule Conditions and Contexts

A **rule** element must contain one or more conditions. A condition is any of the following, which will be described in more detail later:

- a *date/time expression*
- a *node attribute expression*
- a *resource type expression*
- an *operation type expression*
- another **rule** (allowing for complex combinations of conditions)

Each type of condition is allowed only in certain contexts. Although any given context may contain only one **rule** element, that element may contain any number of conditions, including other **rule** elements.

Rules may be used in the following contexts, which also will be described in more detail later:

- a *location constraint*
- a *cluster_property_set* element (within the *crm_config* element)
- an *instance_attributes* element (within an *alert*, *bundle*, *clone*, *group*, *node*, *op*, *primitive*, *recipient*, or *template* element)
- a *meta_attributes* element (within an *alert*, *bundle*, *clone*, *group*, *op*, *op_defaults*, *primitive*, *recipient*, *rsc_defaults*, or *template* element)
- a *utilization* element (within a *node*, *primitive*, or *template* element)

2.11.3 Date/Time Expressions

The *date_expression* element configures a rule condition based on the current date and time. It is allowed in rules in any context.

It may contain a *date_spec* or *duration* element depending on the *operation* as described below.

Table 30: Attributes of a *date_expression* Element

Name	Type	Default	Description
id	<i>id</i>		A unique name for this element (required)
start	<i>ISO 8601</i>		The beginning of the desired time range. Meaningful with an <i>operation</i> of <i>in_range</i> or <i>gt</i> .
end	<i>ISO 8601</i>		The end of the desired time range. Meaningful with an <i>operation</i> of <i>in_range</i> or <i>lt</i> .

Continued on next page

Table 30 – continued from previous page

Name	Type	Default	Description
operation	<i>enumeration</i>	in_range	<p>Specifies how to compare the current date/time against a desired time range. Allowed values:</p> <ul style="list-style-type: none"> • gt: The expression is satisfied if the current date/time is after start (which is required) • lt: The expression is satisfied if the current date/time is before end (which is required) • in_range: The expression is satisfied if the current date/time is greater than or equal to start (if specified) and less than or equal to either end (if specified) or start plus the value of the <i>duration</i> element (if one is contained in the <i>date_expression</i>). At least one of start or end must be specified. If both end and duration are specified, duration is ignored. • date_spec: The expression is satisfied if the current date/time matches the specification given in the contained <i>date_spec</i> element (which is required)

Date Specifications

A *date_spec* element is used within a *date_expression* to specify a combination of dates and times that satisfy the expression.

Table 31: Attributes of a *date_spec* Element

Name	Type	Default	Description
id	<i>id</i>		A unique name for this element (required)
seconds	<i>range</i>		If this is set, the expression is satisfied only if the current time's second is within this range. Allowed integers: 0 to 59.
minutes	<i>range</i>		If this is set, the expression is satisfied only if the current time's minute is within this range. Allowed integers: 0 to 59.
hours	<i>range</i>		If this is set, the expression is satisfied only if the current time's hour is within this range. Allowed integers: 0 to 23 where 0 is midnight and 23 is 11 p.m.
monthdays	<i>range</i>		If this is set, the expression is satisfied only if the current date's day of the month is in this range. Allowed integers: 1 to 31.
weekdays	<i>range</i>		If this is set, the expression is satisfied only if the current date's ordinal day of the week is in this range. Allowed integers: 1-7 (where 1 is Monday and 7 is Sunday).
yeardays	<i>range</i>		If this is set, the expression is satisfied only if the current date's ordinal day of the year is in this range. Allowed integers: 1-366.
months	<i>range</i>		If this is set, the expression is satisfied only if the current date's month is in this range. Allowed integers: 1-12 where 1 is January and 12 is December.

Continued on next page

Table 31 – continued from previous page

Name	Type	Default	Description
weeks	<i>range</i>		If this is set, the expression is satisfied only if the current date's ordinal week of the year is in this range. Allowed integers: 1-53.
years	<i>range</i>		If this is set, the expression is satisfied only if the current date's year according to the Gregorian calendar is in this range.
weekyears	<i>range</i>		If this is set, the expression is satisfied only if the current date's year in which the week started (according to the ISO 8601 standard) is in this range.
moon	<i>range</i>		If this is set, the expression is satisfied only if the current date's phase of the moon is in this range. Allowed values are 0 to 7 where 0 is the new moon and 4 is the full moon. (<i>deprecated since 2.1.6</i>)

Note: Pacemaker can calculate when evaluation of a `date_expression` with an operation of `gt`, `lt`, or `in_range` will next change, and schedule a cluster re-check for that time. However, it does not do this for `date_spec`. Instead, it evaluates the `date_spec` whenever a cluster re-check naturally happens via a cluster event or the `cluster-recheck-interval` cluster option.

For example, if you have a `date_spec` enabling a resource from 9 a.m. to 5 p.m., and `cluster-recheck-interval` has been set to 5 minutes, then sometime between 9 a.m. and 9:05 a.m. the cluster would notice that it needs to start the resource, and sometime between 5 p.m. and 5:05 p.m. it would realize that it needs to stop the resource. The timing of the actual start and stop actions will further depend on factors such as any other actions the cluster may need to perform first, and the load of the machine.

Durations

A `duration` element is used within a `date_expression` to calculate an ending value for `in_range` operations when `end` is not supplied.

Table 32: Attributes of a duration Element

Name	Type	Default	Description
id	<i>id</i>		A unique name for this element (required)
seconds	<i>integer</i>	0	Number of seconds to add to the total duration
minutes	<i>integer</i>	0	Number of minutes to add to the total duration
hours	<i>integer</i>	0	Number of hours to add to the total duration
days	<i>integer</i>	0	Number of days to add to the total duration
weeks	<i>integer</i>	0	Number of weeks to add to the total duration
months	<i>integer</i>	0	Number of months to add to the total duration
years	<i>integer</i>	0	Number of years to add to the total duration

Example Date/Time Expressions

Satisfied if the current year is 2005

```
<rule id="rule1" score="INFINITY">
  <date_expression id="date_expr1" start="2005-001" operation="in_range">
    <duration id="duration1" years="1"/>
  </date_expression>
</rule>
```

or equivalently:

```
<rule id="rule2" score="INFINITY">
  <date_expression id="date_expr2" operation="date_spec">
    <date_spec id="date_spec2" years="2005"/>
  </date_expression>
</rule>
```

9 a.m. to 5 p.m. Monday through Friday

```
<rule id="rule3" score="INFINITY">
  <date_expression id="date_expr3" operation="date_spec">
    <date_spec id="date_spec3" hours="9-16" weekdays="1-5"/>
  </date_expression>
</rule>
```

Note that the 16 matches all the way through 16:59:59, because the numeric value of the hour still matches.

9 a.m. to 6 p.m. Monday through Friday, or anytime Saturday

```
<rule id="rule4" score="INFINITY" boolean-op="or">
  <date_expression id="date_expr4-1" operation="date_spec">
    <date_spec id="date_spec4-1" hours="9-16" weekdays="1-5"/>
  </date_expression>
  <date_expression id="date_expr4-2" operation="date_spec">
    <date_spec id="date_spec4-2" weekdays="6"/>
  </date_expression>
</rule>
```

9 a.m. to 5 p.m. or 9 p.m. to 12 a.m. Monday through Friday

```
<rule id="rule5" score="INFINITY" boolean-op="and">
  <rule id="rule5-nested1" score="INFINITY" boolean-op="or">
    <date_expression id="date_expr5-1" operation="date_spec">
      <date_spec id="date_spec5-1" hours="9-16"/>
    </date_expression>
    <date_expression id="date_expr5-2" operation="date_spec">
      <date_spec id="date_spec5-2" hours="21-23"/>
    </date_expression>
  </rule>
  <date_expression id="date_expr5-3" operation="date_spec">
    <date_spec id="date_spec5-3" weekdays="1-5"/>
  </date_expression>
</rule>
```

Mondays in March 2005

```
<rule id="rule6" score="INFINITY" boolean-op="and">
  <date_expression id="date_expr6-1" operation="date_spec">
    <date_spec id="date_spec6" weekdays="1"/>
  </date_expression>
  <date_expression id="date_expr6-2" operation="in_range"
    start="2005-03-01" end="2005-04-01"/>
</date_expression>
</rule>
```

Note: Because no time is specified with the above dates, 00:00:00 is implied. This means that the range includes all of 2005-03-01 but only the first second of 2005-04-01. You may wish to write `end` as "2005-03-31T23:59:59" to avoid confusion.

2.11.4 Node Attribute Expressions

The `expression` element configures a rule condition based on the value of a node attribute. It is allowed in rules in location constraints and in `instance_attributes` elements within `bundle`, `clone`, `group`, `op`, `primitive`, and `template` elements.

Table 33: Attributes of an expression Element

Name	Type	Default	Description
<code>id</code>	<i>id</i>		A unique name for this element (required)
<code>attribute</code>	<i>text</i>		Name of the node attribute to test (required)

Continued on next page

Table 33 – continued from previous page

Name	Type	Default	Description
operation	<i>enumeration</i>		<p>The comparison to perform (required). Allowed values:</p> <ul style="list-style-type: none"> • defined: The expression is satisfied if the node has the named attribute • not_defined: The expression is satisfied if the node does not have the named attribute • lt: The expression is satisfied if the node attribute value is less than the reference value • gt: The expression is satisfied if the node attribute value is greater than the reference value • lte: The expression is satisfied if the node attribute value is less than or equal to the reference value • gte: The expression is satisfied if the node attribute value is greater than or equal to the reference value • eq: The expression is satisfied if the node attribute value is equal to the reference value • ne: The expression is satisfied if the node attribute value is not equal to the reference value
type	<i>enumeration</i>	The default type for lt , gt , lte , and gte operations is number if either value contains a decimal point character, or integer otherwise. The default type for all other operations is string . If a numeric parse fails for either value, then the values are compared as type string .	How to interpret values. Allowed values are string , integer (<i>since 2.0.5</i>), number , and version . integer truncates floating-point values if necessary before performing a 64-bit integer comparison. number performs a double-precision floating-point comparison (<i>32-bit integer before 2.0.5</i>).
value	<i>text</i>		Reference value to compare node attribute against (used only with, and required for, operations other than defined and not_defined)
value-source	<i>enumeration</i>	literal	<p>How the reference value is obtained. Allowed values:</p> <ul style="list-style-type: none"> • literal: value contains the literal reference value to compare • param: value contains the name of a resource parameter to compare (valid only in the context of a location constraint) • meta: value is the name of a resource meta-attribute to compare (valid only in the context of a location constraint)

In addition to custom node attributes defined by the administrator, the cluster defines special, built-in node attributes for each node that can also be used in rule expressions.

Table 34: Built-in Node Attributes

Name	Description
#uname	<i>Node name</i>
#id	Node ID
#kind	Node type (<code>cluster</code> for cluster nodes, <code>remote</code> for Pacemaker Remote nodes created with the <code>ocf:pacemaker:remote</code> resource, and <code>container</code> for Pacemaker Remote guest nodes and bundle nodes)
#is_dc	<code>true</code> if this node is the cluster's Designated Controller (DC), <code>false</code> otherwise
#cluster-name	The value of the <code>cluster-name</code> cluster property, if set
#site-name	The value of the <code>site-name</code> node attribute, if set, otherwise identical to <code>#cluster-name</code>

2.11.5 Resource Type Expressions

The `rsc_expression` element (*since 2.0.5*) configures a rule condition based on the agent used for a resource. It is allowed in rules in a `meta_attributes` element within a `rsc_defaults` or `op_defaults` element.

Table 35: Attributes of a `rsc_expression` Element

Name	Type	Default	Description
id	<i>id</i>		A unique name for this element (required)
class	<i>text</i>		If this is set, the expression is satisfied only if the resource's agent standard matches this value
provider	<i>text</i>		If this is set, the expression is satisfied only if the resource's agent provider matches this value
type	<i>text</i>		If this is set, the expression is satisfied only if the resource's agent type matches this value

Example Resource Type Expressions

Satisfied for `ocf:heartbeat:IPaddr2` resources

```
<rule id="rule1" score="INFINITY">
  <rsc_expression id="rule_expr1" class="ocf" provider="heartbeat" type="IPaddr2"/>
</rule>
```

Satisfied for `stonith:fence_xvm` resources

```
<rule id="rule2" score="INFINITY">
  <rsc_expression id="rule_expr2" class="stonith" type="fence_xvm"/>
</rule>
```

2.11.6 Operation Type Expressions

The `op_expression` element (*since 2.0.5*) configures a rule condition based on a resource operation name and interval. It is allowed in rules in a `meta_attributes` element within an `op_defaults` element.

Table 36: Attributes of an `op_expression` Element

Name	Type	Default	Description
<code>id</code>	<i>id</i>		A unique name for this element (required)
<code>name</code>	<i>text</i>		The expression is satisfied only if the operation's name matches this value (required)
<code>interval</code>	<i>duration</i>		If this is set, the expression is satisfied only if the operation's interval matches this value

Example Operation Type Expressions

Expression is satisfied for all monitor actions

```
<rule id="rule1" score="INFINITY">
  <op_expression id="rule_expr1" name="monitor"/>
</rule>
```

Expression is satisfied for all monitor actions with a 10-second interval

```
<rule id="rule2" score="INFINITY">
  <op_expression id="rule_expr2" name="monitor" interval="10s"/>
</rule>
```

2.11.7 Using Rules to Determine Resource Location

If a *location constraint* contains a rule, the cluster will apply the constraint to all nodes where the rule is satisfied. This acts as if identical location constraints without rules were defined for each of the nodes.

In the context of a location constraint, `rule` elements may take additional attributes. These have an effect only when set for the constraint's top-level `rule`; they are ignored if set on a subrule.

Table 37: Extra Attributes of a `rule` Element in a Location Constraint

Name	Type	Default	Description
<code>role</code>	<i>enumeration</i>	Started	If this is set in the constraint's top-level rule, the constraint acts as if <code>role</code> were set to this in the <code>rsc_location</code> element.
<code>score</code>	<i>score</i>		If this is set in the constraint's top-level rule, the constraint acts as if <code>score</code> were set to this in the <code>rsc_location</code> element. Only one of <code>score</code> and <code>score-attribute</code> may be set.
<code>score-attribute</code>	<i>text</i>		If this is set in the constraint's top-level rule, the constraint acts as if <code>score</code> were set to the value of this node attribute on each node where the rule is satisfied. Only one of <code>score</code> and <code>score-attribute</code> may be set.

Consider the following simple location constraint:

Prevent resource `webserver` from running on node `node3`

```
<rsc_location id="ban-apache-on-node3" rsc="webserver"
  score="-INFINITY" node="node3"/>
```

The same constraint can be written more verbosely using a rule:

Prevent resource `webserver` from running on node `node3` using a rule

```
<rsc_location id="ban-apache-on-node3" rsc="webserver">
  <rule id="ban-apache-rule" score="-INFINITY">
    <expression id="ban-apache-expr" attribute="#uname"
      operation="eq" value="node3"/>
  </rule>
</rsc_location>
```

The advantage of using the expanded form is that one could add more expressions (for example, limiting the constraint to certain days of the week).

Location Rules Based on Other Node Properties

The expanded form allows us to match node attributes other than its name. As an example, consider this configuration of custom node attributes specifying each node's CPU capacity:

Sample node section with node attributes

```
<nodes>
  <node id="uuid1" uname="c001n01" type="normal">
    <instance_attributes id="uuid1-custom_attrs">
      <nvpair id="uuid1-cpu_mips" name="cpu_mips" value="1234"/>
    </instance_attributes>
  </node>
  <node id="uuid2" uname="c001n02" type="normal">
    <instance_attributes id="uuid2-custom_attrs">
      <nvpair id="uuid2-cpu_mips" name="cpu_mips" value="5678"/>
    </instance_attributes>
  </node>
</nodes>
```

We can use a rule to prevent a resource from running on underpowered machines:

Rule using a node attribute (to be used inside a location constraint)

```
<rule id="need-more-power-rule" score="-INFINITY">
  <expression id="need-more-power-expr" attribute="cpu_mips"
    operation="lt" value="3000"/>
</rule>
```

Using score-attribute Instead of score

When using `score-attribute` instead of `score`, each node matched by the rule has its score adjusted according to its value for the named node attribute.

In the previous example, if the location constraint rule used `score-attribute="cpu_mips"` instead of `score="-INFINITY"`, node `c001n01` would have its preference to run the resource increased by 1234 whereas node `c001n02` would have its preference increased by 5678.

Specifying location scores using pattern submatches

Location constraints may use *rsc-pattern* to apply the constraint to all resources whose IDs match the given pattern. The pattern may contain up to 9 submatches in parentheses, whose values may be used as %1 through %9 in a rule element's `score-attribute` or an expression element's attribute.

For example, the following configuration excerpt gives the resources `server-httpd` and `ip-httpd` a preference of 100 on `node1` and 50 on `node2`, and `ip-gateway` a preference of -100 on `node1` and 200 on `node2`.

Location constraint using submatches

```
<nodes>
  <node id="1" uname="node1">
    <instance_attributes id="node1-attrs">
      <nvpair id="node1-prefer-httpd" name="prefer-httpd" value="100"/>
      <nvpair id="node1-prefer-gateway" name="prefer-gateway" value="-100"/>
    </instance_attributes>
  </node>
  <node id="2" uname="node2">
    <instance_attributes id="node2-attrs">
      <nvpair id="node2-prefer-httpd" name="prefer-httpd" value="50"/>
      <nvpair id="node2-prefer-gateway" name="prefer-gateway" value="200"/>
    </instance_attributes>
  </node>
</nodes>
<resources>
  <primitive id="server-httpd" class="ocf" provider="heartbeat" type="apache"/>
  <primitive id="ip-httpd" class="ocf" provider="heartbeat" type="IPaddr2"/>
  <primitive id="ip-gateway" class="ocf" provider="heartbeat" type="IPaddr2"/>
</resources>
<constraints>
  <!-- The following constraint says that for any resource whose name
  starts with "server-" or "ip-", that resource's preference for a
  node is the value of the node attribute named "prefer-" followed
  by the part of the resource name after "server-" or "ip-",
  wherever such a node attribute is defined.
  -->
  <rsc_location id="location1" rsc-pattern="(server|ip)-(.*)">
    <rule id="location1-rule1" score-attribute="prefer-%2">
      <expression id="location1-rule1-expression1" attribute="prefer-%2" operation="defined"/>
    </rule>
  </rsc_location>
</constraints>
```


2.11.8 Using Rules to Define Options

Rules may be used to control a variety of options:

- *Cluster options* (as `cluster_property_set` elements)
- *Node attributes* (as `instance_attributes` or `utilization` elements inside a `node` element)
- *Resource options* (as `utilization`, `meta_attributes`, or `instance_attributes` elements inside a resource definition element or `op`, `rsc_defaults`, `op_defaults`, or `template` element)
- *Operation options* (as `meta_attributes` elements inside an `op` or `op_defaults` element)
- *Alert options* (as `instance_attributes` or `meta_attributes` elements inside an `alert` or `recipient` element)

Using Rules to Control Resource Options

Often some cluster nodes will be different from their peers. Sometimes, these differences (for example, the location of a binary, or the names of network interfaces) require resources to be configured differently depending on the machine they're hosted on.

By defining multiple `instance_attributes` elements for the resource and adding a rule to each, we can easily handle these special cases.

In the example below, `mySpecialRsc` will use `eth1` and port `9999` when run on `node1`, `eth2` and port `8888` on `node2` and default to `eth0` and port `9999` for all other nodes.

Defining different resource options based on the node name

```
<primitive id="mySpecialRsc" class="ocf" type="Special" provider="me">
  <instance_attributes id="special-node1" score="3">
    <rule id="node1-special-case" score="INFINITY" >
      <expression id="node1-special-case-expr" attribute="#uname"
        operation="eq" value="node1"/>
    </rule>
    <nvpair id="node1-interface" name="interface" value="eth1"/>
  </instance_attributes>
  <instance_attributes id="special-node2" score="2" >
    <rule id="node2-special-case" score="INFINITY">
      <expression id="node2-special-case-expr" attribute="#uname"
        operation="eq" value="node2"/>
    </rule>
    <nvpair id="node2-interface" name="interface" value="eth2"/>
    <nvpair id="node2-port" name="port" value="8888"/>
  </instance_attributes>
  <instance_attributes id="defaults" score="1" >
    <nvpair id="default-interface" name="interface" value="eth0"/>
    <nvpair id="default-port" name="port" value="9999"/>
  </instance_attributes>
</primitive>
```

Multiple `instance_attributes` elements are evaluated from highest score to lowest. If not supplied, the score defaults to zero. Objects with equal scores are processed in their listed order. If an `instance_attributes` object has no rule or a satisfied rule, then for any parameter the resource does not yet have a value for, the resource will use the value defined by the `instance_attributes`.

For example, given the configuration above, if the resource is placed on `node1`:

- `special-node1` has the highest score (3) and so is evaluated first; its rule is satisfied, so `interface` is set to `eth1`.
- `special-node2` is evaluated next with score 2, but its rule is not satisfied, so it is ignored.
- `defaults` is evaluated last with score 1, and has no rule, so its values are examined; `interface` is already defined, so the value here is not used, but `port` is not yet defined, so `port` is set to 9999.

Using Rules to Control Resource Defaults

Rules can be used for resource and operation defaults.

The following example illustrates how to set a different `resource-stickiness` value during and outside work hours. This allows resources to automatically move back to their most preferred hosts, but at a time that (in theory) does not interfere with business activities.

Change `resource-stickiness` during working hours

```
<rsc_defaults>
  <meta_attributes id="core-hours" score="2">
    <rule id="core-hour-rule" score="0">
      <date_expression id="nine-to-five-Mon-to-Fri" operation="date_spec">
        <date_spec id="nine-to-five-Mon-to-Fri-spec" hours="9-16" weekdays="1-5"/>
      </date_expression>
    </rule>
    <nvpair id="core-stickiness" name="resource-stickiness" value="INFINITY"/>
  </meta_attributes>
  <meta_attributes id="after-hours" score="1" >
    <nvpair id="after-stickiness" name="resource-stickiness" value="0"/>
  </meta_attributes>
</rsc_defaults>
```

`rsc_expression` is valid within both `rsc_defaults` and `op_defaults`; `op_expression` is valid only within `op_defaults`.

Default all `IPAddr2` resources to stopped

```
<rsc_defaults>
  <meta_attributes id="op-target-role">
    <rule id="op-target-role-rule" score="INFINITY">
      <rsc_expression id="op-target-role-expr" class="ocf" provider="heartbeat"
        type="IPAddr2"/>
    </rule>
    <nvpair id="op-target-role-nvpair" name="target-role" value="Stopped"/>
  </meta_attributes>
</rsc_defaults>
```

Default all monitor action timeouts to 7 seconds

```

<op_defaults>
  <meta_attributes id="op-monitor-defaults">
    <rule id="op-monitor-default-rule" score="INFINITY">
      <op_expression id="op-monitor-default-expr" name="monitor"/>
    </rule>
    <nvpair id="op-monitor-timeout" name="timeout" value="7s"/>
  </meta_attributes>
</op_defaults>

```

Default the timeout on all 10-second-interval monitor actions on IPAddr2 resources to 8 seconds

```

<op_defaults>
  <meta_attributes id="op-monitor-and">
    <rule id="op-monitor-and-rule" score="INFINITY">
      <rsc_expression id="op-monitor-and-rsc-expr" class="ocf" provider="heartbeat"
        type="IPAddr2"/>
      <op_expression id="op-monitor-and-op-expr" name="monitor" interval="10s"/>
    </rule>
    <nvpair id="op-monitor-and-timeout" name="timeout" value="8s"/>
  </meta_attributes>
</op_defaults>

```

Using Rules to Control Cluster Options

Controlling cluster options is achieved in much the same manner as specifying different resource options on different nodes.

The following example illustrates how to set `maintenance_mode` during a scheduled maintenance window. This will keep the cluster running but not monitor, start, or stop resources during this time.

Schedule a maintenance window for 9 to 11 p.m. CDT Sept. 20, 2019

```

<crm_config>
  <cluster_property_set id="cib-bootstrap-options">
    <nvpair id="bootstrap-stonith-enabled" name="stonith-enabled" value="1"/>
  </cluster_property_set>
  <cluster_property_set id="normal-set" score="10">
    <nvpair id="normal-maintenance-mode" name="maintenance-mode" value="false"/>
  </cluster_property_set>
  <cluster_property_set id="maintenance-window-set" score="1000">
    <nvpair id="maintenance-nvpair1" name="maintenance-mode" value="true"/>
    <rule id="maintenance-rule1" score="INFINITY">
      <date_expression id="maintenance-date1" operation="in_range"
        start="2019-09-20 21:00:00 -05:00" end="2019-09-20 23:00:00 -05:00"/>
    </rule>
  </cluster_property_set>
</crm_config>

```

Important: The `cluster_property_set` with an `id` set to “`cib-bootstrap-options`” will *always* have the

highest priority, regardless of any scores. Therefore, rules in another `cluster_property_set` can never take effect for any properties listed in the bootstrap set.

2.12 Access Control Lists (ACLs)

By default, the `root` user or any user in the `haclient` group can modify Pacemaker's CIB without restriction. Pacemaker offers *access control lists (ACLs)* to provide more fine-grained authorization.

Important: Being able to modify the CIB's resource section allows a user to run any executable file as root, by configuring it as an LSB resource with a full path.

2.12.1 ACL Prerequisites

In order to use ACLs:

- The `enable-acl` *cluster option* must be set to true.
- Desired users must have user accounts in the `haclient` group on all cluster nodes in the cluster.
- If your CIB was created before Pacemaker 1.1.12, it might need to be updated to the current schema (using `cibadmin --upgrade` or a higher-level tool equivalent) in order to use the syntax documented here.
- Prior to the 2.1.0 release, the Pacemaker software had to have been built with ACL support. If you are using an older release, your installation supports ACLs only if the output of the command `pacemakerd --features` contains `acls`. In newer versions, ACLs are always enabled.

Important: `enable-acl` should be set either by the root user, or as part of a batch of CIB changes including roles and users. Otherwise, the user setting it might lock themselves out from making any further changes.

2.12.2 ACL Configuration

ACLs are specified within an `acls` element of the CIB. The `acls` element may contain any number of `acl_role`, `acl_target`, and `acl_group` elements.

2.12.3 ACL Roles

An ACL *role* is a collection of permissions allowing or denying access to particular portions of the CIB. A role is configured with an `acl_role` element in the CIB `acls` section.

Table 38: Properties of an `acl_role` element

Attribute	Description
<code>id</code>	A unique name for the role (<i>required</i>)
<code>description</code>	Arbitrary text for user's use (ignored by Pacemaker)

An `acl_role` element may contain any number of `acl_permission` elements.

Table 39: Properties of an `acl_permission` element

Attribute	Description
<code>id</code>	A unique name for the permission (<i>required</i>)
<code>description</code>	Arbitrary text for user's use (ignored by Pacemaker)
<code>kind</code>	The access being granted. Allowed values are <code>read</code> , <code>write</code> , and <code>deny</code> . A value of <code>write</code> grants both read and write access.
<code>object-type</code>	The name of an XML element in the CIB to which the permission applies. (Exactly one of <code>object-type</code> , <code>xpath</code> , and <code>reference</code> must be specified for a permission.)
<code>attribute</code>	If specified, the permission applies only to <code>object-type</code> elements that have this attribute set (to any value). If not specified, the permission applies to all <code>object-type</code> elements. May only be used with <code>object-type</code> .
<code>reference</code>	The ID of an XML element in the CIB to which the permission applies. (Exactly one of <code>object-type</code> , <code>xpath</code> , and <code>reference</code> must be specified for a permission.)
<code>xpath</code>	An <code>XPath</code> specification selecting an XML element in the CIB to which the permission applies. Attributes may be specified in the <code>XPath</code> to select particular elements, but the permissions apply to the entire element. (Exactly one of <code>object-type</code> , <code>xpath</code> , and <code>reference</code> must be specified for a permission.)

Important:

- Permissions are applied to the selected XML element's entire XML subtree (all elements enclosed within it).
 - Write permission grants the ability to create, modify, or remove the element and its subtree, and also the ability to create any "scaffolding" elements (enclosing elements that do not have attributes other than an ID).
 - Permissions for more specific matches (more deeply nested elements) take precedence over more general ones.
 - If multiple permissions are configured for the same match (for example, in different roles applied to the same user), any `deny` permission takes precedence, then `write`, then lastly `read`.
-

2.12.4 ACL Targets and Groups

ACL targets correspond to user accounts on the system.

Table 40: **Properties of an `acl_target` element**

Attribute	Description
<code>id</code>	A unique identifier for the target (if <code>name</code> is not specified, this must be the name of the user account) (<i>required</i>)
<code>name</code>	If specified, the user account name (this allows you to specify a user name that is already used as the <code>id</code> for some other configuration element) (<i>since 2.1.5</i>)

ACL groups correspond to groups on the system. Any role configured for these groups apply to all users in that group (*since 2.1.5*).

Table 41: **Properties of an `acl_group` element**

Attribute	Description
<code>id</code>	A unique identifier for the group (if <code>name</code> is not specified, this must be the group name) (<i>required</i>)
<code>name</code>	If specified, the group name (this allows you to specify a group name that is already used as the <code>id</code> for some other configuration element)

Each `acl_target` and `acl_group` element may contain any number of `role` elements.

Note: If the system users and groups are defined by some network service (such as LDAP), the cluster itself will be unaffected by outages in the service, but affected users and groups will not be able to make changes to the CIB.

Table 42: **Properties of a `role` element**

Attribute	Description
<code>id</code>	The <code>id</code> of an <code>acl_role</code> element that specifies permissions granted to the enclosing target or group.

Important: The `root` and `hacluster` user accounts always have full access to the CIB, regardless of ACLs. For all other user accounts, when `enable-acl` is true, permission to all parts of the CIB is denied by default (permissions must be explicitly granted).

2.12.5 ACLs and Pacemaker Remote Nodes

ACLs apply differently on Pacemaker Remote nodes, which are assumed to be special-purpose hosts without typical user accounts. Instead, CIB modifications coming from a Pacemaker Remote node use the node's name as the ACL user name, and `pacemaker-remote` as the role.

2.12.6 ACL Examples

```

<acls>

  <acl_role id="read_all">
    <acl_permission id="read_all-cib" kind="read" xpath="/cib" />
  </acl_role>

  <acl_role id="operator">

    <acl_permission id="operator-maintenance-mode" kind="write"
      xpath="//crm_config//nvpair[@name='maintenance-mode']" />

    <acl_permission id="operator-maintenance-attr" kind="write"
      xpath="//nvpair[@name='maintenance']" />

    <acl_permission id="operator-target-role" kind="write"
      xpath="//resources//meta_attributes/nvpair[@name='target-role']" />

    <acl_permission id="operator-is-managed" kind="write"
      xpath="//resources//nvpair[@name='is-managed']" />

    <acl_permission id="operator-rsc_location" kind="write"
      object-type="rsc_location" />

  </acl_role>

  <acl_role id="administrator">
    <acl_permission id="administrator-cib" kind="write" xpath="/cib" />
  </acl_role>

  <acl_role id="minimal">

    <acl_permission id="minimal-standby" kind="read"
      description="allow reading standby node attribute (permanent or transient)"
      xpath="//instance_attributes/nvpair[@name='standby']"/>

    <acl_permission id="minimal-maintenance" kind="read"
      description="allow reading maintenance node attribute (permanent or transient)"
      xpath="//nvpair[@name='maintenance']"/>

    <acl_permission id="minimal-target-role" kind="read"
      description="allow reading resource target roles"
      xpath="//resources//meta_attributes/nvpair[@name='target-role']"/>

    <acl_permission id="minimal-is-managed" kind="read"
      description="allow reading resource managed status"
      xpath="//resources//meta_attributes/nvpair[@name='is-managed']"/>

    <acl_permission id="minimal-deny-instance-attributes" kind="deny"
      xpath="//instance_attributes"/>

    <acl_permission id="minimal-deny-meta-attributes" kind="deny"
      xpath="//meta_attributes"/>

    <acl_permission id="minimal-deny-operations" kind="deny"
      xpath="//operations"/>

    <acl_permission id="minimal-deny-utilization" kind="deny"

```

(continues on next page)

(continued from previous page)

```
    xpath="//utilization"/>

    <acl_permission id="minimal-nodes" kind="read"
      description="allow reading node names/IDs (attributes are denied separately)"
      xpath="/cib/configuration/nodes"/>

    <acl_permission id="minimal-resources" kind="read"
      description="allow reading resource names/agents (parameters are denied separately)"
      xpath="/cib/configuration/resources"/>

    <acl_permission id="minimal-deny-constraints" kind="deny"
      xpath="/cib/configuration/constraints"/>

    <acl_permission id="minimal-deny-topology" kind="deny"
      xpath="/cib/configuration/fencing-topology"/>

    <acl_permission id="minimal-deny-op_defaults" kind="deny"
      xpath="/cib/configuration/op_defaults"/>

    <acl_permission id="minimal-deny-rsc_defaults" kind="deny"
      xpath="/cib/configuration/rsc_defaults"/>

    <acl_permission id="minimal-deny-alerts" kind="deny"
      xpath="/cib/configuration/alerts"/>

    <acl_permission id="minimal-deny-acls" kind="deny"
      xpath="/cib/configuration/acls"/>

    <acl_permission id="minimal-cib" kind="read"
      description="allow reading cib element and crm_config/status sections"
      xpath="/cib"/>

  </acl_role>

  <acl_target id="alice">
    <role id="minimal"/>
  </acl_target>

  <acl_target id="bob">
    <role id="read_all"/>
  </acl_target>

  <acl_target id="carol">
    <role id="read_all"/>
    <role id="operator"/>
  </acl_target>

  <acl_target id="dave">
    <role id="administrator"/>
  </acl_target>

</acls>
```

In the above example, the user `alice` has the minimal permissions necessary to run basic Pacemaker CLI tools, including using `crm_mon` to view the cluster status, without being able to modify anything. The user `bob` can view the entire configuration and status of the cluster, but not make any changes. The user `carol` can read everything, and change selected cluster properties as well as resource roles and location constraints.

Finally, `dave` has full read and write access to the entire CIB.

Looking at the `minimal` role in more depth, it is designed to allow read access to the `cib` tag itself, while denying access to particular portions of its subtree (which is the entire CIB).

This is because the DC node is indicated in the `cib` tag, so `crm_mon` will not be able to report the DC otherwise. However, this does change the security model to allow by default, since any portions of the CIB not explicitly denied will be readable. The `cib` read access could be removed and replaced with read access to just the `crm_config` and `status` sections, for a safer approach at the cost of not seeing the DC in status output.

For a simpler configuration, the `minimal` role allows read access to the entire `crm_config` section, which contains cluster properties. It would be possible to allow read access to specific properties instead (such as `stonith-enabled`, `dc-uuid`, `have-quorum`, and `cluster-name`) to restrict access further while still allowing status output, but cluster properties are unlikely to be considered sensitive.

2.12.7 ACL Limitations

Actions performed via IPC rather than the CIB

ACLs apply *only* to the CIB.

That means ACLs apply to command-line tools that operate by reading or writing the CIB, such as `crm_attribute` when managing permanent node attributes, `crm_mon`, and `cibadmin`.

However, command-line tools that communicate directly with Pacemaker daemons via IPC are not affected by ACLs. For example, users in the `haclient` group may still do the following, regardless of ACLs:

- Query transient node attribute values using `crm_attribute` and `attrd_updater`.
- Query basic node information using `crm_node`.
- Erase resource operation history using `crm_resource`.
- Query fencing configuration information, and execute fencing against nodes, using `stonith_admin`.

ACLs and Pacemaker Remote

ACLs apply to commands run on Pacemaker Remote nodes using the Pacemaker Remote node's name as the ACL user name.

The idea is that Pacemaker Remote nodes (especially virtual machines and containers) are likely to be purpose-built and have different user accounts from full cluster nodes.

2.13 Alerts

Alerts may be configured to take some external action when a cluster event occurs (node failure, resource starting or stopping, etc.).

2.13.1 Alert Agents

As with resource agents, the cluster calls an external program (an *alert agent*) to handle alerts. The cluster passes information about the event to the agent via environment variables. Agents can do anything desired with this information (send an e-mail, log to a file, update a monitoring system, etc.).

Simple alert configuration

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh" />
  </alerts>
</configuration>
```

In the example above, the cluster will call `my-script.sh` for each event.

Multiple alert agents may be configured; the cluster will call all of them for each event.

Alert agents will be called only on cluster nodes. They will be called for events involving Pacemaker Remote nodes, but they will never be called *on* those nodes.

For more information about sample alert agents provided by Pacemaker and about developing custom alert agents, see the *Pacemaker Administration* document.

2.13.2 Alert Recipients

Usually, alerts are directed towards a recipient. Thus, each alert may be additionally configured with one or more recipients. The cluster will call the agent separately for each recipient.

Alert configuration with recipient

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <recipient id="my-alert-recipient" value="some-address"/>
    </alert>
  </alerts>
</configuration>
```

In the above example, the cluster will call `my-script.sh` for each event, passing the recipient `some-address` as an environment variable.

The recipient may be anything the alert agent can recognize – an IP address, an e-mail address, a file name, whatever the particular agent supports.

2.13.3 Alert Meta-Attributes

As with resources, meta-attributes can be configured for alerts to change whether and how Pacemaker calls them.

Table 43: Meta-Attributes of an Alert or Recipient

Meta-Attribute	Default	Description
description		Arbitrary text for user's use (ignored by Pacemaker)
enabled	true	If false for an alert, the alert will not be used. If true for an alert and false for a particular recipient of that alert, that recipient will not be used. (<i>since 2.1.6</i>)

Continued on next page

Table 43 – continued from previous page

Meta-Attribute	Default	Description
timestamp-format	%H:%M:%S.%06N	Format the cluster will use when sending the event's timestamp to the agent. This is a string as used with the <code>date(1)</code> command.
timeout	30s	If the alert agent does not complete within this amount of time, it will be terminated.

Meta-attributes can be configured per alert and/or per recipient.

Alert configuration with meta-attributes

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <meta_attributes id="my-alert-attributes">
        <nvpair id="my-alert-attributes-timeout" name="timeout"
          value="15s"/>
      </meta_attributes>
      <recipient id="my-alert-recipient1" value="someuser@example.com">
        <meta_attributes id="my-alert-recipient1-attributes">
          <nvpair id="my-alert-recipient1-timestamp-format"
            name="timestamp-format" value="%D %H:%M"/>
        </meta_attributes>
      </recipient>
      <recipient id="my-alert-recipient2" value="otheruser@example.com">
        <meta_attributes id="my-alert-recipient2-attributes">
          <nvpair id="my-alert-recipient2-timestamp-format"
            name="timestamp-format" value="%c"/>
        </meta_attributes>
      </recipient>
    </alert>
  </alerts>
</configuration>
```

In the above example, the `my-script.sh` will get called twice for each event, with each call using a 15-second timeout. One call will be passed the recipient `someuser@example.com` and a timestamp in the format `%D %H:%M`, while the other call will be passed the recipient `otheruser@example.com` and a timestamp in the format `%c`.

2.13.4 Alert Instance Attributes

As with resource agents, agent-specific configuration values may be configured as instance attributes. These will be passed to the agent as additional environment variables. The number, names and allowed values of these instance attributes are completely up to the particular agent.

Alert configuration with instance attributes

```

<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <meta_attributes id="my-alert-attributes">
        <nvpair id="my-alert-attributes-timeout" name="timeout"
          value="15s"/>
      </meta_attributes>
      <instance_attributes id="my-alert-options">
        <nvpair id="my-alert-options-debug" name="debug"
          value="false"/>
      </instance_attributes>
      <recipient id="my-alert-recipient1"
        value="someuser@example.com"/>
    </alert>
  </alerts>
</configuration>

```

2.13.5 Alert Filters

By default, an alert agent will be called for node events, fencing events, and resource events. An agent may choose to ignore certain types of events, but there is still the overhead of calling it for those events. To eliminate that overhead, you may select which types of events the agent should receive.

Alert filters are configured within a `select` element inside an `alert` element.

Table 44: Possible alert filters

Name	Events alerted
<code>select_nodes</code>	A node joins or leaves the cluster (whether at the cluster layer for cluster nodes, or via a remote connection for Pacemaker Remote nodes).
<code>select_fencing</code>	Fencing or unfencing of a node completes (whether successfully or not).
<code>select_resources</code>	A resource action other than meta-data completes (whether successfully or not).
<code>select_attributes</code>	A transient attribute value update is sent to the CIB.

Alert configuration to receive only node events and fencing events

```

<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <select>
        <select_nodes />
        <select_fencing />
      </select>
      <recipient id="my-alert-recipient1"
        value="someuser@example.com"/>
    </alert>
  </alerts>
</configuration>

```

With `<select_attributes>` (the only event type not enabled by default), the agent will receive alerts when a node attribute changes. If you wish the agent to be called only when certain attributes change, you can configure that as well.

Alert configuration to be called when certain node attributes change

```

<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <select>
        <select_attributes>
          <attribute id="alert-standby" name="standby" />
          <attribute id="alert-shutdown" name="shutdown" />
        </select_attributes>
      </select>
      <recipient id="my-alert-recipient1" value="someuser@example.com"/>
    </alert>
  </alerts>
</configuration>

```

Node attribute alerts are currently considered experimental. Alerts may be limited to attributes set via `attrd_updater`, and agents may be called multiple times with the same attribute value.

2.14 Reusing Parts of the Configuration

Pacemaker provides multiple ways to simplify the configuration XML by reusing parts of it in multiple places.

Besides simplifying the XML, this also allows you to manipulate multiple configuration elements with a single reference.

2.14.1 Reusing Resource Definitions

If you want to create lots of resources with similar configurations, defining a *resource template* simplifies the task. Once defined, it can be referenced in primitives or in certain types of constraints.

Configuring Resources with Templates

The primitives referencing the template will inherit all meta-attributes, instance attributes, utilization attributes and operations defined in the template. And you can define specific attributes and operations for any of the primitives. If any of these are defined in both the template and the primitive, the values defined in the primitive will take precedence over the ones defined in the template.

Hence, resource templates help to reduce the amount of configuration work. If any changes are needed, they can be done to the template definition and will take effect globally in all resource definitions referencing that template.

Resource templates have a syntax similar to that of primitives.

Resource template for a migratable Xen virtual machine

```
<template id="vm-template" class="ocf" provider="heartbeat" type="Xen">
  <meta_attributes id="vm-template-meta_attributes">
    <nvpair id="vm-template-meta_attributes-allow-migrate" name="allow-migrate" value="true"/>
  </meta_attributes>
  <utilization id="vm-template-utilization">
    <nvpair id="vm-template-utilization-memory" name="memory" value="512"/>
  </utilization>
  <operations>
    <op id="vm-template-monitor-15s" interval="15s" name="monitor" timeout="60s"/>
    <op id="vm-template-start-0" interval="0" name="start" timeout="60s"/>
  </operations>
</template>
```

Once you define a resource template, you can use it in primitives by specifying the `template` property.

Xen primitive resource using a resource template

```
<primitive id="vm1" template="vm-template">
  <instance_attributes id="vm1-instance_attributes">
    <nvpair id="vm1-instance_attributes-name" name="name" value="vm1"/>
    <nvpair id="vm1-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shared-vm/vm1"/>
  </instance_attributes>
</primitive>
```

In the example above, the new primitive `vm1` will inherit everything from `vm-template`. For example, the equivalent of the above two examples would be:

Equivalent Xen primitive resource not using a resource template

```
<primitive id="vm1" class="ocf" provider="heartbeat" type="Xen">
  <meta_attributes id="vm-template-meta_attributes">
    <nvpair id="vm-template-meta_attributes-allow-migrate" name="allow-migrate" value="true"/>
  </meta_attributes>
  <utilization id="vm-template-utilization">
    <nvpair id="vm-template-utilization-memory" name="memory" value="512"/>
  </utilization>
  <operations>
    <op id="vm-template-monitor-15s" interval="15s" name="monitor" timeout="60s"/>
    <op id="vm-template-start-0" interval="0" name="start" timeout="60s"/>
  </operations>
  <instance_attributes id="vm1-instance_attributes">
    <nvpair id="vm1-instance_attributes-name" name="name" value="vm1"/>
    <nvpair id="vm1-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shared-vm/vm1"/>
  </instance_attributes>
</primitive>
```

If you want to overwrite some attributes or operations, add them to the particular primitive's definition.

Xen resource overriding template values

```
<primitive id="vm2" template="vm-template">
  <meta_attributes id="vm2-meta_attributes">
    <nvpair id="vm2-meta_attributes-allow-migrate" name="allow-migrate" value="false"/>
  </meta_attributes>
  <utilization id="vm2-utilization">
    <nvpair id="vm2-utilization-memory" name="memory" value="1024"/>
  </utilization>
  <instance_attributes id="vm2-instance_attributes">
    <nvpair id="vm2-instance_attributes-name" name="name" value="vm2"/>
    <nvpair id="vm2-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shared-vm/vm2"/>
  </instance_attributes>
  <operations>
    <op id="vm2-monitor-30s" interval="30s" name="monitor" timeout="120s"/>
    <op id="vm2-stop-0" interval="0" name="stop" timeout="60s"/>
  </operations>
</primitive>
```

In the example above, the new primitive `vm2` has special attribute values. Its `monitor` operation has a longer timeout and interval, and the primitive has an additional `stop` operation.

To see the resulting definition of a resource, run:

```
# crm_resource --query-xml --resource vm2
```

To see the raw definition of a resource in the CIB, run:

```
# crm_resource --query-xml-raw --resource vm2
```

Using Templates in Constraints

A resource template can be referenced in the following types of constraints:

- `order` constraints (see *Specifying the Order in which Resources Should Start/Stop*)
- `colocation` constraints (see *Placing Resources Relative to other Resources*)
- `rsc_ticket` constraints (for multi-site clusters as described in *Configuring Ticket Dependencies*)

Resource templates referenced in constraints stand for all primitives which are derived from that template. This means, the constraint applies to all primitive resources referencing the resource template. Referencing resource templates in constraints is an alternative to resource sets and can simplify the cluster configuration considerably.

For example, given the example templates earlier in this chapter:

```
<rsc_colocation id="vm-template-colo-base-rsc" rsc="vm-template" rsc-role="Started" with-rsc="base-
→rsc" score="INFINITY"/>
```

would collocate all VMs with `base-rsc` and is the equivalent of the following constraint configuration:

```
<rsc_colocation id="vm-colo-base-rsc" score="INFINITY">
  <resource_set id="vm-colo-base-rsc-0" sequential="false" role="Started">
    <resource_ref id="vm1"/>
    <resource_ref id="vm2"/>
  </resource_set>
  <resource_set id="vm-colo-base-rsc-1">
```

(continues on next page)

(continued from previous page)

```

    <resource_ref id="base-rsc"/>
  </resource_set>
</rsc_colocation>

```

Note: In a colocation constraint, only one template may be referenced from either `rsc` or `with-rsc`; the other reference must be a regular resource.

Using Templates in Resource Sets

Resource templates can also be referenced in resource sets.

For example, given the example templates earlier in this section, then:

```

<rsc_order id="order1" score="INFINITY">
  <resource_set id="order1-0">
    <resource_ref id="base-rsc"/>
    <resource_ref id="vm-template"/>
    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>

```

is the equivalent of the following constraint using a sequential resource set:

```

<rsc_order id="order1" score="INFINITY">
  <resource_set id="order1-0">
    <resource_ref id="base-rsc"/>
    <resource_ref id="vm1"/>
    <resource_ref id="vm2"/>
    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>

```

Or, if the resources referencing the template can run in parallel, then:

```

<rsc_order id="order2" score="INFINITY">
  <resource_set id="order2-0">
    <resource_ref id="base-rsc"/>
  </resource_set>
  <resource_set id="order2-1" sequential="false">
    <resource_ref id="vm-template"/>
  </resource_set>
  <resource_set id="order2-2">
    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>

```

is the equivalent of the following constraint configuration:

```

<rsc_order id="order2" score="INFINITY">
  <resource_set id="order2-0">
    <resource_ref id="base-rsc"/>
  </resource_set>
  <resource_set id="order2-1" sequential="false">

```

(continues on next page)

(continued from previous page)

```

<resource_ref id="vm1"/>
<resource_ref id="vm2"/>
</resource_set>
<resource_set id="order2-2">
  <resource_ref id="top-rsc"/>
</resource_set>
</rsc_order>

```

2.14.2 Reusing Rules, Options and Sets of Operations

Sometimes a number of constraints need to use the same set of rules, and resources need to set the same options and parameters. To simplify this situation, you can refer to an existing object using an `id-ref` instead of an `id`.

So if for one resource you have

```

<rsc_location id="WebServer-connectivity" rsc="Webserver">
  <rule id="ping-prefer-rule" score-attribute="pingd" >
    <expression id="ping-prefer" attribute="pingd" operation="defined"/>
  </rule>
</rsc_location>

```

Then instead of duplicating the rule for all your other resources, you can instead specify:

Referencing rules from other constraints

```

<rsc_location id="WebDB-connectivity" rsc="WebDB">
  <rule id-ref="ping-prefer-rule"/>
</rsc_location>

```

Important: The cluster will insist that the `rule` exists somewhere. Attempting to add a reference to a nonexistent `id` will cause a validation failure, as will attempting to remove a `rule` with an `id` that is referenced elsewhere.

Some rule syntax is allowed only in *certain contexts*. Validation cannot ensure that the referenced rule is allowed in the context of the rule containing `id-ref`, so such errors will be caught (and logged) only after the new configuration is accepted. It is the administrator's responsibility to check for these.

The same principle applies for `meta_attributes` and `instance_attributes` as illustrated in the example below:

Referencing attributes, options, and operations from other resources

```

<primitive id="mySpecialRsc" class="ocf" type="Special" provider="me">
  <instance_attributes id="mySpecialRsc-attrs" score="1" >
    <nvpair id="default-interface" name="interface" value="eth0"/>
    <nvpair id="default-port" name="port" value="9999"/>
  </instance_attributes>
  <meta_attributes id="mySpecialRsc-options">
    <nvpair id="failure-timeout" name="failure-timeout" value="5m"/>
    <nvpair id="migration-threshold" name="migration-threshold" value="1"/>
    <nvpair id="stickiness" name="resource-stickiness" value="0"/>
  </meta_attributes>
  <operations id="health-checks">
    <op id="health-check" name="monitor" interval="60s"/>
    <op id="health-check" name="monitor" interval="30min"/>
  </operations>
</primitive>
<primitive id="myOtherRsc" class="ocf" type="Other" provider="me">
  <instance_attributes id-ref="mySpecialRsc-attrs"/>
  <meta_attributes id-ref="mySpecialRsc-options"/>
  <operations id-ref="health-checks"/>
</primitive>

```

`id-ref` can similarly be used with `resource_set` (in any constraint type), `nvpair`, and `operations`.

2.14.3 Tagging Configuration Elements

Pacemaker allows you to *tag* any configuration element that has an XML ID.

The main purpose of tagging is to support higher-level user interface tools; Pacemaker itself only uses tags within constraints. Therefore, what you can do with tags mostly depends on the tools you use.

Configuring Tags

A tag is simply a named list of XML IDs.

Tag referencing three resources

```

<tags>
  <tag id="all-vms">
    <obj_ref id="vm1"/>
    <obj_ref id="vm2"/>
    <obj_ref id="vm3"/>
  </tag>
</tags>

```

What you can do with this new tag depends on what your higher-level tools support. For example, a tool might allow you to enable or disable all of the tagged resources at once, or show the status of just the tagged resources.

A single configuration element can be listed in any number of tags.

Important: If listing nodes in a tag, you must list the node's `id`, not `name`.

Using Tags in Constraints and Resource Sets

Pacemaker itself only uses tags in constraints. If you supply a tag name instead of a resource name in any constraint, the constraint will apply to all resources listed in that tag.

Constraint using a tag

```
<rsc_order id="order1" first="storage" then="all-vms" kind="Mandatory" />
```

In the example above, assuming the `all-vms` tag is defined as in the previous example, the constraint will behave the same as:

Equivalent constraints without tags

```
<rsc_order id="order1-1" first="storage" then="vm1" kind="Mandatory" />
<rsc_order id="order1-2" first="storage" then="vm2" kind="Mandatory" />
<rsc_order id="order1-3" first="storage" then="vm3" kind="Mandatory" />
```

A tag may be used directly in the constraint, or indirectly by being listed in a *resource set* used in the constraint. When used in a resource set, an expanded tag will honor the set's `sequential` property.

Filtering With Tags

The `crm_mon` tool can be used to display lots of information about the state of the cluster. On large or complicated clusters, this can include a lot of information, which makes it difficult to find the one thing you are interested in. The `--resource=` and `--node=` command line options can be used to filter results. In their most basic usage, these options take a single resource or node name. However, they can also be supplied with a tag name to display several objects at once.

For instance, given the following CIB section:

```
<resources>
  <primitive class="stonith" id="Fencing" type="fence_xvm"/>
  <primitive class="ocf" id="dummy" provider="pacemaker" type="Dummy"/>
  <group id="inactive-group">
    <primitive class="ocf" id="inactive-dummy-1" provider="pacemaker" type="Dummy"/>
    <primitive class="ocf" id="inactive-dummy-2" provider="pacemaker" type="Dummy"/>
  </group>
  <clone id="inactive-clone">
    <primitive id="inactive-dhcpd" class="systemd" type="dhcpd"/>
  </clone>
</resources>
<tags>
  <tag id="inactive-rscs">
    <obj_ref id="inactive-group"/>
    <obj_ref id="inactive-clone"/>
  </tag>
</tags>
```

The following would be output for `crm_mon --resource=inactive-rscs -r:`

```

Cluster Summary:
* Stack: corosync
* Current DC: cluster02 (version 2.0.4-1.e97f9675f.git.e17-e97f9675f) - partition with quorum
* Last updated: Tue Oct 20 16:09:01 2020
* Last change: Tue May 5 12:04:36 2020 by hacluster via crmd on cluster01
* 5 nodes configured
* 27 resource instances configured (4 DISABLED)

Node List:
* Online: [ cluster01 cluster02 ]

Full List of Resources:
* Clone Set: inactive-clone [inactive-dhcpd] (disabled):
  * Stopped (disabled): [ cluster01 cluster02 ]
* Resource Group: inactive-group (disabled):
  * inactive-dummy-1 (ocf::pacemaker:Dummy): Stopped (disabled)
  * inactive-dummy-2 (ocf::pacemaker:Dummy): Stopped (disabled)

```

2.15 Status

Pacemaker automatically generates a `status` section in the CIB (inside the `cib` element, at the same level as `configuration`). The status is transient, and is not stored to disk with the rest of the CIB.

The section's structure and contents are internal to Pacemaker and subject to change from release to release. Its often obscure element and attribute names are kept for historical reasons, to maintain compatibility with older versions during rolling upgrades.

Users should not modify the section directly, though various command-line tool options affect it indirectly.

2.15.1 Node State

The `status` element contains `node_state` elements for each node in the cluster (and potentially nodes that have been removed from the configuration since the cluster started). The `node_state` element has attributes that allow the cluster to determine whether the node is healthy.

Example minimal node state entry

```

<node_state id="1" uname="cl-virt-1" in_ccm="1721760952" crmd="1721760952" crm-debug-origin=
↳"controld_update_resource_history" join="member" expected="member">
  <transient_attributes id="1"/>
  <lrmd id="1"/>
</node_state>

```

Table 45: Attributes of a `node_state` Element

Name	Type	Description
id	<i>text</i>	Node ID (identical to <code>id</code> of corresponding <code>node</code> element in the <code>configuration</code> section)
uname	<i>text</i>	Node name (identical to <code>uname</code> of corresponding <code>node</code> element in the <code>configuration</code> section)

Continued on next page

Table 45 – continued from previous page

Name	Type	Description
in_ccm	<i>epoch time (since 2.1.7; previously boolean)</i>	If the node's controller is currently in the cluster layer's membership, this is the epoch time at which it joined (or 1 if the node is in the process of leaving the cluster), otherwise 0 (<i>since 2.1.7; previously, it was "true" or "false"</i>)
crmd	<i>epoch time (since 2.1.7; previously an enumeration)</i>	If the node's controller is currently in the cluster layer's controller messaging group, this is the epoch time at which it joined, otherwise 0 (<i>since 2.1.7; previously, the value was either "online" or "offline"</i>)
crm-debug-origin	<i>text</i>	Name of the source code function that recorded this <code>node_state</code> element (for debugging)
join	<i>enumeration</i>	Current status of node's controller join sequence (and thus whether it is eligible to run resources). Allowed values: <ul style="list-style-type: none"> • <code>down</code>: Not yet joined • <code>pending</code>: In the process of joining or leaving • <code>member</code>: Fully joined • <code>banned</code>: Rejected by DC
expected	<i>enumeration</i>	What cluster expects <code>join</code> to be in the immediate future. Allowed values are same as for <code>join</code> .

2.15.2 Transient Node Attributes

The `transient_attributes` section specifies transient *Node Attributes*. In addition to any values set by the administrator or resource agents using the `atrd_updater` or `crm_attribute` tools, the cluster stores various state information here.

Example transient node attributes for a node

```
<transient_attributes id="cl-virt-1">
  <instance_attributes id="status-cl-virt-1">
    <nvpair id="status-cl-virt-1-pingd" name="pingd" value="3"/>
    <nvpair id="status-cl-virt-1-fail-count-pingd:0.monitor_30000" name="fail-count-pingd:0
↪#monitor_30000" value="1"/>
    <nvpair id="status-cl-virt-1-last-failure-pingd:0" name="last-failure-pingd:0" value=
↪"1239009742"/>
  </instance_attributes>
</transient_attributes>
```

2.15.3 Node History

Each `node_state` element contains an `lrm` element with a history of certain resource actions performed on the node. The `lrm` element contains an `lrm_resources` element.

Resource History

The `lrm_resources` element contains an `lrm_resource` element for each resource that has had an action performed on the node.

An `lrm_resource` entry has attributes allowing the cluster to stop the resource safely even if it is removed from the configuration. Specifically, the resource's `id`, `class`, `type` and `provider` are recorded.

Action History

Each `lrm_resource` element contains an `lrm_rsc_op` element for each recorded action performed for that resource on that node. (Not all actions are recorded, just enough to determine the resource's state.)

Table 46: Attributes of an `lrm_rsc_op` element

Name	Type	Description
<code>id</code>	<i>text</i>	Identifier for the history entry constructed from the resource ID, action name or history entry type, and action interval.
<code>operation_key</code>	<i>text</i>	Identifier for the action that was executed, constructed from the resource ID, action name, and action interval.
<code>operation</code>	<i>text</i>	The name of the action the history entry is for
<code>crm-debug-origin</code>	<i>text</i>	Name of the source code function that recorded this entry (for debugging)
<code>crm_feature_set</code>	<i>version</i>	The Pacemaker feature set used to record this entry.
<code>transition-key</code>	<i>text</i>	A concatenation of the action's transition graph action number, the transition graph number, the action's expected result, and the UUID of the controller instance that scheduled it.
<code>transition-magic</code>	<i>text</i>	A concatenation of <code>op-status</code> , <code>rc-code</code> , and <code>transition-key</code> .
<code>exit-reason</code>	<i>text</i>	An error message (if available) from the resource agent or Pacemaker if the action did not return success.
<code>on_node</code>	<i>text</i>	The name of the node that executed the action (identical to the <code>uname</code> of the enclosing <code>node_state</code> element)
<code>call-id</code>	<i>integer</i>	A node-specific counter used to determine the order in which actions were executed.
<code>rc-code</code>	<i>integer</i>	The resource agent's exit status for this action. Refer to the <i>Resource Agents</i> chapter of <i>Pacemaker Administration</i> for how these values are interpreted.
<code>op-status</code>	<i>integer</i>	The execution status of this action. The meanings of these codes are internal to Pacemaker.
<code>interval</code>	<i>nonnegative integer</i>	If the action is recurring, its frequency (in milliseconds), otherwise 0.
<code>last-rc-change</code>	<i>epoch time</i>	Node-local time at which the action first returned the current value of <code>rc-code</code> .
<code>exec-time</code>	<i>integer</i>	Time (in seconds) that action execution took (if known)
<code>queue-time</code>	<i>integer</i>	Time (in seconds) that action was queued in the local executor (if known)
<code>op-digest</code>	<i>text</i>	If present, this is a hash of the parameters passed to the action. If a hash of the currently configured parameters does not match this, that means the resource configuration changed since the action was performed, and the resource must be reloaded or restarted.
<code>op-restart-digest</code>	<i>text</i>	If present, the resource agent supports reloadable parameters, and this is a hash of the non-reloadable parameters passed to the action. This allows the cluster to choose between reload and restart when one is needed.

Continued on next page

Table 46 – continued from previous page

Name	Type	Description
op-secure-digest	<i>text</i>	If present, the resource agent marks some parameters as sensitive, and this is a hash of the non-sensitive parameters passed to the action. This allows the value of sensitive parameters to be removed from a saved copy of the CIB while still allowing scheduler simulations to be performed on that copy.

Simple Operation History Example

A monitor operation (determines current state of the apcstonith resource)

```
<lrms_resource id="apcstonith" type="fence_apc_snmp" class="stonith">
  <lrms_rsc_op id="apcstonith_monitor_0" operation="monitor" call-id="2"
    rc-code="7" op-status="0" interval="0"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
    op-digest="2e3da9274d3550dc6526fb24bfcba0"
    transition-key="22:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    transition-magic="0:7;22:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239008085" exec-time="10" queue-time="0"/>
</lrms_resource>
```

The above example shows the history entry for a probe (non-recurring monitor operation) for the `apcstonith` resource.

The cluster schedules probes for every configured resource on a node when the node first starts, in order to determine the resource's current state before it takes any further action.

From the `transition-key`, we can see that this was the 22nd action of the 2nd graph produced by this instance of the controller (2668bbeb-06d5-40f9-936d-24cb7f87006a).

The third field of the `transition-key` contains a 7, which indicates that the cluster expects to find the resource inactive. By looking at the `rc-code` property, we see that this was the case.

As that is the only action recorded for this node, we can conclude that the cluster started the resource elsewhere.

Complex Operation History Example

Resource history of a pingd clone with multiple entries

```
<lrms_resource id="pingd:0" type="pingd" class="ocf" provider="pacemaker">
  <lrms_rsc_op id="pingd:0_monitor_30000" operation="monitor" call-id="34"
    rc-code="0" op-status="0" interval="30000"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
    transition-key="10:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239009741" exec-time="10" queue-time="0"/>
  <lrms_rsc_op id="pingd:0_stop_0" operation="stop"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1" call-id="32"
    rc-code="0" op-status="0" interval="0"
    transition-key="11:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239009741" exec-time="10" queue-time="0"/>
  <lrms_rsc_op id="pingd:0_start_0" operation="start" call-id="33"
    rc-code="0" op-status="0" interval="0"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
    transition-key="31:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239009741" exec-time="10" queue-time="0" />
  <lrms_rsc_op id="pingd:0_monitor_0" operation="monitor" call-id="3"
    rc-code="0" op-status="0" interval="0"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
```

When more than one history entry exists, it is important to first sort them by `call-id` before interpreting them.

Once sorted, the above example can be summarized as:

1. A non-recurring monitor operation returning 7 (not running), with a `call-id` of 3
2. A stop operation returning 0 (success), with a `call-id` of 32
3. A start operation returning 0 (success), with a `call-id` of 33
4. A recurring monitor returning 0 (success), with a `call-id` of 34

The cluster processes each history entry to build up a picture of the resource's state. After the first and second entries, it is considered stopped, and after the third it considered active.

Based on the last operation, we can tell that the resource is currently active.

Additionally, from the presence of a `stop` operation with a lower `call-id` than that of the `start` operation, we can conclude that the resource has been restarted. Specifically this occurred as part of actions 11 and 31 of transition 11 from the controller instance with the key `2668bbeb...`. This information can be helpful for locating the relevant section of the logs when looking for the source of a failure.

2.16 Multi-Site Clusters and Tickets

Apart from local clusters, Pacemaker also supports multi-site clusters. That means you can have multiple, geographically dispersed sites, each with a local cluster. Failover between these clusters can be coordinated manually by the administrator, or automatically by a higher-level entity called a *Cluster Ticket Registry (CTR)*.

2.16.1 Challenges for Multi-Site Clusters

Typically, multi-site environments are too far apart to support synchronous communication and data replication between the sites. That leads to significant challenges:

- How do we make sure that a cluster site is up and running?
- How do we make sure that resources are only started once?
- How do we make sure that quorum can be reached between the different sites and a split-brain scenario avoided?
- How do we manage failover between sites?
- How do we deal with high latency in case of resources that need to be stopped?

In the following sections, learn how to meet these challenges.

2.16.2 Conceptual Overview

Multi-site clusters can be considered as “overlay” clusters where each cluster site corresponds to a cluster node in a traditional cluster. The overlay cluster can be managed by a CTR in order to guarantee that any cluster resource will be active on no more than one cluster site. This is achieved by using *tickets* that are treated as failover domain between cluster sites, in case a site should be down.

The following sections explain the individual components and mechanisms that were introduced for multi-site clusters in more detail.

Ticket

Tickets are, essentially, cluster-wide attributes. A ticket grants the right to run certain resources on a specific cluster site. Resources can be bound to a certain ticket by `rsc_ticket` constraints. Only if the ticket is available at a site can the respective resources be started there. Vice versa, if the ticket is revoked, the resources depending on that ticket must be stopped.

The ticket thus is similar to a *site quorum*, i.e. the permission to manage/own resources associated with that site. (One can also think of the current `have-quorum` flag as a special, cluster-wide ticket that is granted in case of node majority.)

Tickets can be granted and revoked either manually by administrators (which could be the default for classic enterprise clusters), or via the automated CTR mechanism described below.

A ticket can only be owned by one site at a time. Initially, none of the sites has a ticket. Each ticket must be granted once by the cluster administrator.

The presence or absence of tickets for a site is stored in the CIB as a cluster status. With regards to a certain ticket, there are only two states for a site: `true` (the site has the ticket) or `false` (the site does not have the ticket). The absence of a certain ticket (during the initial state of the multi-site cluster) is the same as the value `false`.

Dead Man Dependency

A site can only activate resources safely if it can be sure that the other site has deactivated them. However after a ticket is revoked, it can take a long time until all resources depending on that ticket are stopped “cleanly”, especially in case of cascaded resources. To cut that process short, the concept of a *Dead Man Dependency* was introduced.

If a dead man dependency is in force, if a ticket is revoked from a site, the nodes that are hosting dependent resources are fenced. This considerably speeds up the recovery process of the cluster and makes sure that resources can be migrated more quickly.

This can be configured by specifying a `loss-policy="fence"` in `rsc_ticket` constraints.

Cluster Ticket Registry

A CTR is a coordinated group of network daemons that automatically handles granting, revoking, and timing out tickets (instead of the administrator revoking the ticket somewhere, waiting for everything to stop, and then granting it on the desired site).

Pacemaker does not implement its own CTR, but interoperates with external software designed for that purpose (similar to how resource and fencing agents are not directly part of pacemaker).

Participating clusters run the CTR daemons, which connect to each other, exchange information about their connectivity, and vote on which sites gets which tickets.

A ticket is granted to a site only once the CTR is sure that the ticket has been relinquished by the previous owner, implemented via a timer in most scenarios. If a site loses connection to its peers, its tickets time out and recovery occurs. After the connection timeout plus the recovery timeout has passed, the other sites are allowed to re-acquire the ticket and start the resources again.

This can also be thought of as a “quorum server”, except that it is not a single quorum ticket, but several.

Configuration Replication

As usual, the CIB is synchronized within each cluster, but it is *not* synchronized across cluster sites of a multi-site cluster. You have to configure the resources that will be highly available across the multi-site cluster for every site accordingly.

2.16.3 Configuring Ticket Dependencies

The `rsc_ticket` constraint lets you specify the resources depending on a certain ticket. Together with the constraint, you can set a **loss-policy** that defines what should happen to the respective resources if the ticket is revoked.

The attribute **loss-policy** can have the following values:

- **fence**: Fence the nodes that are running the relevant resources.
- **stop**: Stop the relevant resources.
- **freeze**: Do nothing to the relevant resources.
- **demote**: Demote relevant resources that are running in the promoted role.

Constraint that fences node if ticketA is revoked

```
<rsc_ticket id="rsc1-req-ticketA" rsc="rsc1" ticket="ticketA" loss-policy="fence"/>
```

The example above creates a constraint with the ID `rsc1-req-ticketA`. It defines that the resource `rsc1` depends on `ticketA` and that the node running the resource should be fenced if `ticketA` is revoked.

If resource `rsc1` were a promotable resource, you might want to configure that only being in the promoted role depends on `ticketA`. With the following configuration, `rsc1` will be demoted if `ticketA` is revoked:

Constraint that demotes rsc1 if ticketA is revoked

```
<rsc_ticket id="rsc1-req-ticketA" rsc="rsc1" rsc-role="Promoted" ticket="ticketA" loss-policy="demote"/>
```

You can create multiple `rsc_ticket` constraints to let multiple resources depend on the same ticket. However, `rsc_ticket` also supports resource sets (see *Resource Sets*), so one can easily list all the resources in one `rsc_ticket` constraint instead.

Ticket constraint for multiple resources

```
<rsc_ticket id="resources-dep-ticketA" ticket="ticketA" loss-policy="fence">
  <resource_set id="resources-dep-ticketA-0" role="Started">
    <resource_ref id="rsc1"/>
    <resource_ref id="group1"/>
    <resource_ref id="clone1"/>
  </resource_set>
  <resource_set id="resources-dep-ticketA-1" role="Promoted">
    <resource_ref id="ms1"/>
  </resource_set>
</rsc_ticket>
```

In the example above, there are two resource sets, so we can list resources with different roles in a single `rsc_ticket` constraint. There's no dependency between the two resource sets, and there's no dependency among the resources within a resource set. Each of the resources just depends on `ticketA`.

Referencing resource templates in `rsc_ticket` constraints, and even referencing them within resource sets, is also supported.

If you want other resources to depend on further tickets, create as many constraints as necessary with `rsc_ticket`.

2.16.4 Managing Multi-Site Clusters

Granting and Revoking Tickets Manually

You can grant tickets to sites or revoke them from sites manually. If you want to re-distribute a ticket, you should wait for the dependent resources to stop cleanly at the previous site before you grant the ticket to the new site.

Use the `crm_ticket` command line tool to grant and revoke tickets.

To grant a ticket to this site:

```
# crm_ticket --ticket ticketA --grant
```

To revoke a ticket from this site:

```
# crm_ticket --ticket ticketA --revoke
```

Important: If you are managing tickets manually, use the `crm_ticket` command with great care, because it cannot check whether the same ticket is already granted elsewhere.

Granting and Revoking Tickets via a Cluster Ticket Registry

We will use `Booth` here as an example of software that can be used with pacemaker as a Cluster Ticket Registry. `Booth` implements the `Raft` algorithm to guarantee the distributed consensus among different cluster sites, and manages the ticket distribution (and thus the failover process between sites).

Each of the participating clusters and *arbitrators* runs the `Booth` daemon `boothd`.

An *arbitrator* is the multi-site equivalent of a quorum-only node in a local cluster. If you have a setup with an even number of sites, you need an additional instance to reach consensus about decisions such as failover of resources across sites. In this case, add one or more arbitrators running at additional sites. Arbitrators are single machines that run a `booth` instance in a special mode. An arbitrator is especially important for a two-site scenario, otherwise there is no way for one site to distinguish between a network failure between it and the other site, and a failure of the other site.

The most common multi-site scenario is probably a multi-site cluster with two sites and a single arbitrator on a third site. However, technically, there are no limitations with regards to the number of sites and the number of arbitrators involved.

`Boothd` at each site connects to its peers running at the other sites and exchanges connectivity details. Once a ticket is granted to a site, the `booth` mechanism will manage the ticket automatically: If the site which holds the ticket is out of service, the `booth` daemons will vote which of the other sites will get the ticket. To protect against brief connection failures, sites that lose the vote (either explicitly or implicitly by being disconnected from the voting body) need to relinquish the ticket after a time-out. Thus, it is made

sure that a ticket will only be re-distributed after it has been relinquished by the previous site. The resources that depend on that ticket will fail over to the new site holding the ticket. The nodes that have run the resources before will be treated according to the **loss-policy** you set within the **rsc_ticket** constraint.

Before the booth can manage a certain ticket within the multi-site cluster, you initially need to grant it to a site manually via the **booth** command-line tool. After you have initially granted a ticket to a site, **boothd** will take over and manage the ticket automatically.

Important: The **booth** command-line tool can be used to grant, list, or revoke tickets and can be run on any machine where **boothd** is running. If you are managing tickets via Booth, use only **booth** for manual intervention, not **crm_ticket**. That ensures the same ticket will only be owned by one cluster site at a time.

Booth Requirements

- All clusters that will be part of the multi-site cluster must be based on Pacemaker.
- Booth must be installed on all cluster nodes and on all arbitrators that will be part of the multi-site cluster.
- Nodes belonging to the same cluster site should be synchronized via NTP. However, time synchronization is not required between the individual cluster sites.

General Management of Tickets

Display the information of tickets:

```
# crm_ticket --info
```

Or you can monitor them with:

```
# crm_mon --tickets
```

Display the **rsc_ticket** constraints that apply to a ticket:

```
# crm_ticket --ticket ticketA --constraints
```

When you want to do maintenance or manual switch-over of a ticket, revoking the ticket would trigger the loss policies. If **loss-policy="fence"**, the dependent resources could not be gracefully stopped/demoted, and other unrelated resources could even be affected.

The proper way is making the ticket *standby* first with:

```
# crm_ticket --ticket ticketA --standby
```

Then the dependent resources will be stopped or demoted gracefully without triggering the loss policies.

If you have finished the maintenance and want to activate the ticket again, you can run:

```
# crm_ticket --ticket ticketA --activate
```

2.16.5 For more information

- SUSE's Geo Clustering quick start
- Booth

2.17 Sample Configurations

2.17.1 Empty

An Empty Configuration

```
<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="0" num_updates=
↪"0">
  <configuration>
    <crm_config/>
    <nodes/>
    <resources/>
    <constraints/>
  </configuration>
  <status/>
</cib>
```

2.17.2 Simple

A simple configuration with two nodes, some cluster options and a resource

```
<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="0" num_updates=
↪"0">
  <configuration>
    <crm_config>
      <cluster_property_set id="cib-bootstrap-options">
        <nvpair id="option-1" name="symmetric-cluster" value="true"/>
        <nvpair id="option-2" name="no-quorum-policy" value="stop"/>
        <nvpair id="option-3" name="stonith-enabled" value="0"/>
      </cluster_property_set>
    </crm_config>
    <nodes>
      <node id="xxx" uname="c001n01" type="normal"/>
      <node id="yyy" uname="c001n02" type="normal"/>
    </nodes>
    <resources>
      <primitive id="myAddr" class="ocf" provider="heartbeat" type="IPaddr">
        <operations>
          <op id="myAddr-monitor" name="monitor" interval="300s"/>
        </operations>
        <instance_attributes id="myAddr-params">
          <nvpair id="myAddr-ip" name="ip" value="192.0.2.10"/>
        </instance_attributes>
      </primitive>
    </resources>
    <constraints>
      <rsc_location id="myAddr-prefer" rsc="myAddr" node="c001n01" score="INFINITY"/>
    </constraints>
    <rsc_defaults>
      <meta_attributes id="rsc_defaults-options">
        <nvpair id="rsc-default-1" name="resource-stickiness" value="100"/>
        <nvpair id="rsc-default-2" name="migration-threshold" value="10"/>
      </meta_attributes>
    </rsc_defaults>
  </configuration>
  <status/>
</cib>
```

In the above example, we have one resource (an IP address) that we check every five minutes and will run on host c001n01 until either the resource fails 10 times or the host shuts down.

2.17.3 Advanced Configuration

An advanced configuration with groups, clones and STONITH

```

<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="0" num_updates=
↪"0">
  <configuration>
    <crm_config>
      <cluster_property_set id="cib-bootstrap-options">
        <nvpair id="option-1" name="symmetric-cluster" value="true"/>
        <nvpair id="option-2" name="no-quorum-policy" value="stop"/>
        <nvpair id="option-3" name="stonith-enabled" value="true"/>
      </cluster_property_set>
    </crm_config>
    <nodes>
      <node id="xxx" uname="c001n01" type="normal"/>
      <node id="yyy" uname="c001n02" type="normal"/>
      <node id="zzz" uname="c001n03" type="normal"/>
    </nodes>
    <resources>
      <primitive id="myAddr" class="ocf" provider="heartbeat" type="IPAddr">
        <operations>
          <op id="myAddr-monitor" name="monitor" interval="300s"/>
        </operations>
        <instance_attributes id="myAddr-attrs">
          <nvpair id="myAddr-attr-1" name="ip" value="192.0.2.10"/>
        </instance_attributes>
      </primitive>
      <group id="myGroup">
        <primitive id="database" class="systemd" type="mariadb">
          <operations>
            <op id="database-monitor" name="monitor" interval="300s"/>
          </operations>
        </primitive>
        <primitive id="webserver" class="systemd" type="httpd">
          <operations>
            <op id="webserver-monitor" name="monitor" interval="300s"/>
          </operations>
        </primitive>
      </group>
      <clone id="STONITH">
        <meta_attributes id="stonith-options">
          <nvpair id="stonith-option-1" name="globally-unique" value="false"/>
        </meta_attributes>
        <primitive id="stonithclone" class="stonith" type="external/ssh">
          <operations>
            <op id="stonith-op-mon" name="monitor" interval="5s"/>
          </operations>
          <instance_attributes id="stonith-attrs">
            <nvpair id="stonith-attr-1" name="hostlist" value="c001n01,c001n02"/>
          </instance_attributes>
        </primitive>
      </clone>
    </resources>
    <constraints>
      <rscl_location id="myAddr-prefer" rsc="myAddr" node="c001n01"
score="INFINITY"/>
      <rscl_colocation id="group-with-ip" rsc="myGroup" with-rsc="myAddr"

```



INDEX

- genindex
- search

Symbols

#digests

node attribute, 34

#node-unfenced

node attribute, 34

A

Access Control List (ACL), 118

acl_group, 120

acl_permission, 119

acl_role, 118

acl_target, 119

acls, 118

role, 120

acl_group

id (attribute), 120

name (attribute), 120

XML element, 120

acl_permission

attribute (attribute), 119

description (attribute), 119, 124

id (attribute), 119

kind (attribute), 119

object-type (attribute), 119

reference (attribute), 119

XML element, 119

xpath (attribute), 119

acl_role

description (attribute), 118

id (attribute), 118

XML element, 118

acl_target

id (attribute), 120

name (attribute), 120

XML element, 119

acls

XML element, 118

action

history, 136

property, description, 46

property, enabled, 47

property, id, 46

property, interval, 46

property, interval-origin, 48

property, name, 46

property, on-fail, 47

property, record-pending, 48

property, role, 46

property, start-delay, 48

property, timeout, 46

resource_set attribute, 63

add-host

network attribute, 97

admin_epoch

cib, 21

agent

alert, 123

alert, 123

agent, 123

filters, 126

instance attributes, 125

meta-attribute, enabled, 124

meta-attribute, timeout, 125

meta-attribute, timestamp-format, 125

meta-attributes, 124

recipient, 124

XML element, 123

alerts

XML element, 123

allow-migrate

resource option, 41

allow-unhealthy-nodes

resource option, 41

Asymmetrical Clusters, 56

attribute

acl_permission attribute, 119

action (resource_set), 63

add-host (network), 97

attribute (acl_permission), 119

control-port (network), 97

description (acl_permission), 119, 124

description (acl_role), 118

description (bundle), 96

description (clone), 90

- description (group), 88
- expression, 109
- first (rsc_order), 58
- first-action (rsc_order), 58
- host-interface (network), 97
- host-netmask (network), 97
- id (acl_group), 120
- id (acl_permission), 119
- id (acl_role), 118
- id (acl_target), 120
- id (bundle), 96
- id (cluster_property_set), 20
- id (instance_attributes), 20
- id (meta_attributes), 20
- id (port-mapping), 98
- id (resource_set), 63
- id (role), 120
- id (rsc_colocation), 60
- id (rsc_location), 55
- id (rsc_order), 58
- id (storage-mapping), 98
- id (utilization), 20
- image (docker), 96
- image (podman), 96
- influence (rsc_colocation), 61
- internal-port (port-mapping), 98
- ip-range-start (network), 97
- kind (acl_permission), 119
- kind (resource_set), 63
- kind (rsc_order), 59
- name (acl_group), 120
- name (acl_target), 120
- network (docker), 96
- network (podman), 96
- node (rsc_location), 55
- node-attribute (rsc_colocation), 60
- object-type (acl_permission), 119
- options (docker), 96
- options (podman), 96
- options (storage-mapping), 98
- port (port-mapping), 98
- promoted-max (docker), 96
- promoted-max (podman), 96
- range (port-mapping), 98
- reference (acl_permission), 119
- replicas (docker), 96
- replicas (podman), 96
- replicas-per-host (docker), 96
- replicas-per-host (podman), 96
- require-all (resource_set), 63
- resource-discovery (rsc_location), 56
- role (resource_set), 63
- role (rsc_location), 56
- rsc (rsc_colocation), 60

- rsc (rsc_location), 55
- rsc-pattern (rsc_location), 55
- run-command (docker), 96
- run-command (podman), 96
- score (cluster_property_set), 20
- score (instance_attributes), 20
- score (meta_attributes), 20
- score (resource_set), 63
- score (rsc_colocation), 60
- score (rsc_location), 55
- score (utilization), 20
- sequential (resource_set), 63
- source-dir (storage-mapping), 98
- source-dir-root (storage-mapping), 98
- symmetrical (rsc_order), 59
- target-dir (storage-mapping), 98
- then (rsc_order), 58
- then-action (rsc_order), 58
- with-rsc (rsc_colocation), 60
- XML element, 126
- xpath (acl_permission), 119

B

- batch-limit
 - cluster option, 23
- boolean
 - type, 10
- boolean-op
 - rule, 105
- bundle
 - attribute, description, 96
 - attribute, id, 96
 - meta-attributes, 100
 - network, 97
 - node attributes, 100
 - primitive, 99
 - XML element, 96

C

- call-id
 - lrm_rsc_op, 136
- cib
 - admin_epoch, 21
 - cib-last-written, 21
 - dc-uuid, 21
 - epoch, 21
 - execution-date, 22
 - have-quorum, 21
 - num_updates, 21
 - remote-clear-port, 21
 - remote-tls-port, 21
 - validate-with, 21
 - XML element, 19
- cib-last-written

- cib, 21
- CIB_pam_service
 - node option, 12
- class
 - resource, 38
 - rsc_expression, 111
- clone, 89
 - attribute, description, 90
 - constraint, 92
 - option, clone-max, 90
 - option, clone-min, 90
 - option, clone-node-max, 90
 - option, globally-unique, 90
 - option, interleave, 91
 - option, notify, 90
 - option, ordered, 91
 - option, promotable, 91
 - option, promoted-max, 91
 - option, promoted-node-max, 91
 - options, 90
 - ordering constraint, rsc-role, 60
 - ordering constraint, with-rsc-role, 60
 - property, id, 90
 - resource-stickiness, 94
 - XML element, 90
- clone-max
 - clone option, 90
- clone-min
 - clone option, 90
- clone-node-max
 - clone option, 90
- cluster option
 - batch-limit, 23
 - cluster-delay, 26
 - cluster-infrastructure, 22
 - cluster-ipc-limit, 26
 - cluster-name, 22
 - cluster-recheck-interval, 28
 - concurrent-fencing, 25
 - dc-deadtime, 26
 - dc-version, 22
 - election-timeout, 29
 - enable-acl, 27
 - enable-startup-probes, 24
 - fence-reaction, 26
 - have-watchdog, 25
 - join-finalization-timeout, 29
 - join-integration-timeout, 29
 - load-threshold, 23
 - maintenance-mode, 24
 - migration-limit, 23
 - no-quorum-policy, 23
 - node-action-limit, 23
 - node-health-base, 27
 - node-health-green, 27
 - node-health-red, 27
 - node-health-strategy, 27, 34
 - node-health-yellow, 27
 - node-pending-timeout, 26
 - pe-error-series-max, 27
 - pe-input-series-max, 27
 - pe-warn-series-max, 27
 - placement-strategy, 27
 - priority-fencing-delay, 26
 - rule, 114, 117
 - shutdown-escalation, 29
 - shutdown-lock, 28
 - shutdown-lock-limit, 28
 - start-failure-is-fatal, 24
 - startup-fencing, 29
 - stonith-action, 24
 - stonith-enabled, 24
 - stonith-max-attempts, 24
 - stonith-timeout, 24
 - stonith-watchdog-timeout, 25
 - stop-all-resources, 23
 - stop-orphan-actions, 24
 - stop-orphan-resources, 23
 - symmetric-cluster, 23
 - transition-delay, 29
- cluster-delay
 - cluster option, 26
- cluster-infrastructure
 - cluster option, 22
- cluster-ipc-limit
 - cluster option, 26
- cluster-name
 - cluster option, 22
- cluster-recheck-interval
 - cluster option, 28
- cluster_property_set
 - id, 20
 - score, 20
- colocation, 59
- concurrent-fencing
 - cluster option, 25
- configuration
 - XML element, 10, 19
- constraint, 54
 - colocation, 59
 - location, 55
 - ordering, 58
 - resource set, 62
 - rsc_colocation, 60
 - rsc_location, 55
 - rsc_order, 58
- container-attribute-target
 - resource option, 41

- control-port
 - network attribute, 97
- critical
 - resource option, 39
- crm-debug-origin
 - lrm_rsc_op, 136
 - node_state, 135
- crm_feature_set
 - lrm_rsc_op, 136
- crmd
 - node_state, 135
- custom
 - node-health-strategy value, 35

D

- date specification, 106
- date/time
 - type, 10
- date_expression
 - end, 105
 - id, 105
 - operation, 106
 - start, 105
 - XML element, 105
- date_spec
 - hours, 106
 - id, 106
 - minutes, 106
 - monthdays, 106
 - months, 106
 - moon, 107
 - seconds, 106
 - weekdays, 106
 - weeks, 107
 - weekyears, 107
 - XML element, 106
 - yeardays, 106
 - years, 107

- days
 - duration, 107
- dc-deadtime
 - cluster option, 26
- dc-uuid
 - cib, 21
- dc-version
 - cluster option, 22
- description
 - acl_permission attribute, 119, 124
 - acl_role attribute, 118
 - action property, 46
 - bundle attribute, 96
 - clone attribute, 90
 - group attribute, 88
 - op, 46

- resource, 38
- devices
 - fencing-level, 84
- docker
 - attribute, image, 96
 - attribute, network, 96
 - attribute, options, 96
 - attribute, promoted-max, 96
 - attribute, replicas, 96
 - attribute, replicas-per-host, 96
 - attribute, run-command, 96
 - XML element, 96
- duration, 107
 - days, 107
 - hours, 107
 - id, 107
 - minutes, 107
 - months, 107
 - seconds, 107
 - type, 10
 - weeks, 107
 - XML element, 107
 - years, 107

E

- election-timeout
 - cluster option, 29
- enable-acl
 - cluster option, 27
- enable-startup-probes
 - cluster option, 24
- enabled
 - action property, 47
 - alert meta-attribute, 124
 - op, 47
- end
 - date_expression, 105
- enumeration
 - type, 11
- epoch
 - cib, 21
- epoch_time
 - type, 11
- exec-time
 - lrm_rsc_op, 136
- execution-date
 - cib, 22
- exit-reason
 - lrm_rsc_op, 136
- expected
 - node_state, 135
- expression
 - attribute, 109
 - id, 109

- operation, 110
 - type, 110
 - value, 110
 - value-source, 110
 - XML element, 109
- ## F
- fail-count
 - node attribute, 33
 - failure-timeout
 - resource option, 41
 - fence-reaction
 - cluster option, 26
 - fencing, 70
 - agent, 71
 - alert, 123
 - configuration, 77
 - device, 70
 - special instance attributes, 71
 - topology, 84
 - unfencing, 76
 - why necessary, 70
 - fencing-level, 84
 - devices, 84
 - id, 84
 - index, 84
 - target, 84
 - target-attribute, 84
 - target-pattern, 84
 - target-value, 84
 - fencing-topology, 84
 - first
 - rsc_order attribute, 58
 - first-action
 - rsc_order attribute, 58
- ## G
- globally-unique
 - clone option, 90
 - green
 - node health attribute value, 34
 - group
 - attribute, description, 88
 - property, id, 88
 - resource-stickiness, 89
 - XML element, 88
 - guest node, 44
- ## H
- have-quorum
 - cib, 21
 - have-watchdog
 - cluster option, 25
 - history
 - action, 136
 - node, 135
 - resource, 135
 - host-interface
 - network attribute, 97
 - host-netmask
 - network attribute, 97
 - hours
 - date_spec, 106
 - duration, 107
- ## I
- id
 - acl_group attribute, 120
 - acl_permission attribute, 119
 - acl_role attribute, 118
 - acl_target attribute, 120
 - action property, 46
 - bundle attribute, 96
 - clone property, 90
 - cluster_property_set, 20
 - date_expression, 105
 - date_spec, 106
 - duration, 107
 - expression, 109
 - fencing-level, 84
 - group property, 88
 - instance_attributes, 20
 - lrm_rsc_op, 136
 - meta_attributes, 20
 - node_state, 134
 - op, 46
 - op_expression, 112
 - port-mapping attribute, 98
 - resource, 38
 - resource_set attribute, 63
 - role attribute, 120
 - rsc_colocation attribute, 60
 - rsc_expression, 111
 - rsc_location attribute, 55
 - rsc_order attribute, 58
 - rule, 104
 - storage-mapping attribute, 98
 - type, 11
 - utilization, 20
 - image
 - docker attribute, 96
 - podman attribute, 96
 - in_ccm
 - node_state, 135
 - index
 - fencing-level, 84
 - influence
 - rsc_colocation attribute, 61

- instance attribute
 - alert instance attributes, 125
 - rule, 114
- instance_attributes
 - id, 20
 - score, 20
- integer
 - type, 11
- interleave
 - clone option, 91
- internal-port
 - port-mapping attribute, 98
- interval
 - action property, 46
 - lrm_rsc_op, 136
 - op, 46
 - op_expression, 112
- interval-origin
 - action property, 48
 - op, 48
- ip-range-start
 - network attribute, 97
- is-managed
 - resource option, 39
- iso8601
 - type, 11

J

- join
 - node_state, 135
- join-finalization-timeout
 - cluster option, 29
- join-integration-timeout
 - cluster option, 29

K

- kind
 - acl_permission attribute, 119
 - resource_set attribute, 63
 - rsc_order attribute, 59

L

- last-failure
 - node attribute, 33
- last-rc-change
 - lrm_rsc_op, 136
- Linux Standard Base
 - resources, 37
- load-threshold
 - cluster option, 23
- location constraint, 55
 - rule, 112
- lrm
 - XML element, 135

- lrm_resource
 - XML element, 135
- lrm_resources
 - XML element, 135
- lrm_rsc_op
 - call-id, 136
 - crm-debug-origin, 136
 - crm_feature_set, 136
 - exec-time, 136
 - exit-reason, 136
 - id, 136
 - interval, 136
 - last-rc-change, 136
 - on_node, 136
 - op-digest, 136
 - op-restart-digest, 136
 - op-secure-digest, 137
 - op-status, 136
 - operation, 136
 - operation_key, 136
 - queue-time, 136
 - rc-code, 136
 - transition-key, 136
 - transition-magic, 136
 - XML element, 136

LSB

- resources, 37

M

- maintenance
 - node attribute, 33
 - resource option, 40
- maintenance-mode
 - cluster option, 24
- meta-attribute
 - alert meta-attributes, 124
 - enabled (alert), 124
 - rule, 114
 - timeout (alert), 125
 - timestamp-format (alert), 125
- meta_attributes
 - id, 20
 - score, 20
- migrate-on-red
 - node-health-strategy value, 35
- migration-limit
 - cluster option, 23
- migration-threshold
 - resource meta-attribute, 52
 - resource option, 40
- minutes
 - date_spec, 106
 - duration, 107
- monthdays

- date_spec, 106
- months
 - date_spec, 106
 - duration, 107
- moon
 - date_spec, 107
- multiple-active
 - resource option, 41
- N**
- name
 - acl_group attribute, 120
 - acl_target attribute, 120
 - action property, 46
 - op, 46
 - op_expression, 112
- network
 - attribute
 - control-port, 97
 - host-interface, 97
 - host-netmask, 97
 - attribute, add-host, 97
 - attribute, ip-range-start, 97
 - docker attribute, 96
 - podman attribute, 96
 - XML element, 97
- no-quorum-policy
 - cluster option, 23
- node, 29
 - alert, 123
 - attribute, 31
 - cluster node, 29
 - guest, 44
 - health, 34
 - history, 135
 - name, 31
 - Pacemaker Remote, 30
 - quorum-only, 31
 - remote, 44
 - rsc_location attribute, 55
 - state, 134
 - transient attribute, 135
- node attribute, 31
 - #digests, 34
 - #node-unfenced, 34
 - fail-count, 33
 - health, 34
 - health (green), 34
 - health (red), 34
 - health (score), 34
 - health (yellow), 34
 - last-failure, 33
 - maintenance, 33
 - probe_complete, 33
 - resource-discovery-enabled, 33
 - rule, 114
 - rule expression, 109
 - shutdown, 33
 - site-name, 33
 - standby, 33
 - terminate, 33
 - transient, 135
- node option
 - CIB_pam_service, 12
 - PCMK_authkey_location, 16
 - PCMK_blackbox, 14
 - PCMK_ca_file, 15
 - PCMK_callgrind_enabled, 18
 - PCMK_cert_file, 15
 - PCMK_cluster_type, 18
 - PCMK_crl_file, 16
 - PCMK_debug, 13
 - PCMK_dh_max_bits, 17
 - PCMK_fail_fast, 14
 - PCMK_ipc_buffer, 18
 - PCMK_ipc_type, 18
 - PCMK_key_file, 16
 - PCMK_logfacility, 12
 - PCMK_logfile, 13
 - PCMK_logfile_mode, 13
 - PCMK_logpriority, 12
 - PCMK_node_action_limit, 14
 - PCMK_node_start_state, 14
 - PCMK_panic_action, 15
 - PCMK_remote_address, 15
 - PCMK_remote_pid1, 17
 - PCMK_remote_port, 15
 - PCMK_remote_schema_directory, 18
 - PCMK_schema_directory, 18
 - PCMK_stderr, 13
 - PCMK_tls_priorities, 17
 - PCMK_trace_blackbox, 14
 - PCMK_trace_files, 13
 - PCMK_trace_formats, 14
 - PCMK_trace_functions, 13
 - PCMK_trace_tags, 14
 - PCMK_valgrind_enabled, 18
 - SBD_SYNC_RESOURCE_STARTUP, 18
 - SBD_WATCHDOG_TIMEOUT, 18
 - VALGRIND_OPTS, 19
- node-action-limit
 - cluster option, 23
- node-attribute
 - rsc_colocation attribute, 60
- node-health-base
 - cluster option, 27
- node-health-green
 - cluster option, 27

- node-health-red
 - cluster option, 27
- node-health-strategy
 - cluster option, 27, 34
 - custom, 35
 - migrate-on-red, 35
 - none, 35
 - only-green, 35
 - progressive, 35
- node-health-yellow
 - cluster option, 27
- node-pending-timeout
 - cluster option, 26
- node_state
 - crm-debug-origin, 135
 - crmd, 135
 - expected, 135
 - id, 134
 - in_ccm, 135
 - join, 135
 - uname, 134
 - XML element, 134
- none
 - node-health-strategy value, 35
- nonnegative integer
 - type, 11
- notify
 - clone option, 90
- num_updates
 - cib, 21

O

- object-type
 - acl_permission attribute, 119
- OCF
 - resources, 36
- on-fail
 - action property, 47
 - op, 47
- on_node
 - lrm_rsc_op, 136
- only-green
 - node-health-strategy value, 35
- op
 - description, 46
 - enabled, 47
 - id, 46
 - interval, 46
 - interval-origin, 48
 - name, 46
 - on-fail, 47
 - record-pending, 48
 - role, 46
 - start-delay, 48

- timeout, 46
- op-digest
 - lrm_rsc_op, 136
- op-restart-digest
 - lrm_rsc_op, 136
- op-secure-digest
 - lrm_rsc_op, 137
- op-status
 - lrm_rsc_op, 136
- op_expression
 - id, 112
 - interval, 112
 - name, 112
 - XML element, 111
- Open Cluster Framework
 - resources, 36
- operation
 - date_expression, 106
 - expression, 110
 - failure count, 52
 - failure recovery, 51
 - lrm_rsc_op, 136
 - rule expression, 111
- operation defaults
 - rule, 114
- operation_key
 - lrm_rsc_op, 136
- Opt-In Clusters, 56
- Opt-Out Clusters, 57
- option
 - clone-max (clone), 90
 - clone-min (clone), 90
 - clone-node-max (clone), 90
 - globally-unique (clone), 90
 - interleave (clone), 91
 - notify (clone), 90
 - ordered (clone), 91
 - promotable (clone), 91
 - promoted-max (clone), 91
 - promoted-node-max (clone), 91
- options
 - clone, 90
 - docker attribute, 96
 - podman attribute, 96
 - rule, 104
 - storage-mapping attribute, 98
- ordered
 - clone option, 91
- ordering constraint
 - rsc-role (clone), 60
 - with-rsc-role (clone), 60

P

- Pacemaker Remote

- guest node, 44
- node, 30
- remote node, 44
- pcmk_action_limit, 73
- PCMK_authkey_location
 - node option, 16
- PCMK_blackbox
 - node option, 14
- PCMK_ca_file
 - node option, 15
- PCMK_callgrind_enabled
 - node option, 18
- PCMK_cert_file
 - node option, 15
- PCMK_cluster_type
 - node option, 18
- PCMK_crl_file
 - node option, 16
- PCMK_debug
 - node option, 13
- pcmk_delay_base, 73
- pcmk_delay_max, 73
- PCMK_dh_max_bits
 - node option, 17
- PCMK_fail_fast
 - node option, 14
- pcmk_host_argument, 73
- pcmk_host_check, 72
- pcmk_host_list, 72
- pcmk_host_map, 72
- PCMK_ipc_buffer
 - node option, 18
- PCMK_ipc_type
 - node option, 18
- PCMK_key_file
 - node option, 16
- pcmk_list_action, 74
- pcmk_list_retries, 74
- pcmk_list_timeout, 74
- PCMK_logfacility
 - node option, 12
- PCMK_logfile
 - node option, 13
- PCMK_logfile_mode
 - node option, 13
- PCMK_logpriority
 - node option, 12
- pcmk_monitor_action, 75
- pcmk_monitor_retries, 75
- pcmk_monitor_timeout, 75
- PCMK_node_action_limit
 - node option, 14
- PCMK_node_start_state
 - node option, 14
- pcmk_off_action, 74
- pcmk_off_retries, 74
- pcmk_off_timeout, 74
- PCMK_panic_action
 - node option, 15
- pcmk_reboot_action, 73
- pcmk_reboot_retries, 74
- pcmk_reboot_timeout, 73
- PCMK_remote_address
 - node option, 15
- PCMK_remote_pid1
 - node option, 17
- PCMK_remote_port
 - node option, 15
- PCMK_remote_schema_directory
 - node option, 18
- PCMK_schema_directory
 - node option, 18
- pcmk_status_action, 75
- pcmk_status_retries, 75
- pcmk_status_timeout, 75
- PCMK_stderr
 - node option, 13
- PCMK_tls_priorities
 - node option, 17
- PCMK_trace_blackbox
 - node option, 14
- PCMK_trace_files
 - node option, 13
- PCMK_trace_formats
 - node option, 14
- PCMK_trace_functions
 - node option, 13
- PCMK_trace_tags
 - node option, 14
- PCMK_valgrind_enabled
 - node option, 18
- pe-error-series-max
 - cluster option, 27
- pe-input-series-max
 - cluster option, 27
- pe-warn-series-max
 - cluster option, 27
- percentage
 - type, 11
- placement-strategy
 - cluster option, 27
- podman
 - attribute, image, 96
 - attribute, network, 96
 - attribute, options, 96
 - attribute, promoted-max, 96
 - attribute, replicas, 96
 - attribute, replicas-per-host, 96

- attribute, run-command, 96
- XML element, 96
- port
 - port-mapping attribute, 98
 - remote node, 44
 - type, 11
- port-mapping
 - attribute, id, 98
 - attribute, internal-port, 98
 - attribute, port, 98
 - attribute, range, 98
 - XML element, 97
- priority
 - resource option, 39
- priority-fencing-delay
 - cluster option, 26
- probe_complete
 - node attribute, 33
- progressive
 - node-health-strategy value, 35
- promotable
 - clone option, 91
- promotable clone, 90
 - constraint, 92
- promoted-max
 - clone option, 91
 - docker attribute, 96
 - podman attribute, 96
- promoted-node-max
 - clone option, 91
- property
 - id (clone), 90
 - id (group), 88
- provider
 - resource, 38
 - rsc_expression, 111
- provides, 72

Q

- queue-time
 - lrm_rsc_op, 136
- quorum-only node, 31

R

- range
 - port-mapping attribute, 98
 - type, 11
- rc-code
 - lrm_rsc_op, 136
- recipient
 - XML element, 124
- reconnect_interval
 - remote node, 44
- record-pending

- action property, 48
 - op, 48
- red
 - node health attribute value, 34
- reference
 - acl_permission attribute, 119
- reload, 53
- reload-agent, 53
- remote node, 44
 - port, 44
 - reconnect_interval, 44
 - server, 44
- remote-addr
 - resource option, 45
- remote-allow-migrate
 - resource option, 45
- remote-clear-port
 - cib, 21
- remote-connect-timeout
 - resource option, 45
- remote-node
 - resource option, 44
- remote-port
 - resource option, 45
- remote-tls-port
 - cib, 21
- replicas
 - docker attribute, 96
 - podman attribute, 96
- replicas-per-host
 - docker attribute, 96
 - podman attribute, 96
- require-all
 - resource_set attribute, 63
- requires
 - resource option, 40
- Resource
 - STONITH, 38
 - System Services, 37
 - Systemd, 37
- resource, 36
 - action, 45
 - alert, 123
 - clone, 89
 - constraint, 54
 - failure count, 52
 - failure recovery, 51
 - history, 135
 - location relative to other resources, 59
 - LSB, 37
 - migration-threshold, 52
 - OCF, 36
 - operation, 45
 - option, allow-migrate, 41

- option, allow-unhealthy-nodes, 41
- option, container-attribute-target, 41
- option, critical, 39
- option, failure-timeout, 41
- option, is-managed, 39
- option, maintenance, 40
- option, migration-threshold, 40
- option, multiple-active, 41
- option, priority, 39
- option, remote-addr, 45
- option, remote-allow-migrate, 45
- option, remote-connect-timeout, 45
- option, remote-node, 44
- option, remote-port, 45
- option, requires, 40
- option, resource-stickiness, 40
- option, target-role, 39
- promotable, 90
- property, class, 38
- property, description, 38
- property, id, 38
- property, provider, 38
- property, type, 38
- resource set, 62
- rule expression, 111
- standard, 36
- start order, 58
- resource defaults
 - rule, 114
- resource-discovery
 - rsc_location attribute, 56
- resource-discovery-enabled
 - node attribute, 33
- resource-stickiness
 - clone, 94
 - group, 89
 - resource option, 40
- resource_set
 - attribute, action, 63
 - attribute, id, 63
 - attribute, kind, 63
 - attribute, require-all, 63
 - attribute, role, 63
 - attribute, score, 63
 - attribute, sequential, 63
 - XML element, 62
- role
 - action property, 46
 - id (attribute), 120
 - op, 46
 - resource_set attribute, 63
 - rsc_location attribute, 56
 - rule, 112
 - XML element, 120
- rsc
 - rsc_colocation attribute, 60
 - rsc_location attribute, 55
- rsc-pattern
 - rsc_location attribute, 55
- rsc-role
 - clone ordering constraint, 60
- rsc_colocation
 - attribute, id, 60
 - attribute, influence, 61
 - attribute, node-attribute, 60
 - attribute, rsc, 60
 - attribute, score, 60
 - attribute, with-rsc, 60
 - XML element, 60
- rsc_expression
 - class, 111
 - id, 111
 - provider, 111
 - type, 111
 - XML element, 111
- rsc_location
 - attribute, id, 55
 - attribute, node, 55
 - attribute, resource-discovery, 56
 - attribute, role, 56
 - attribute, rsc, 55
 - attribute, rsc-pattern, 55
 - attribute, score, 55
 - XML element, 55
- rsc_order
 - attribute, first, 58
 - attribute, first-action, 58
 - attribute, id, 58
 - attribute, kind, 59
 - attribute, symmetrical, 59
 - attribute, then, 58
 - attribute, then-action, 58
 - constraint, 58
 - XML element, 58
- rule, 104
 - boolean-op, 105
 - cluster option, 114, 117
 - conditions, 105
 - contexts, 105
 - date/time expression, 105
 - id, 104
 - instance attribute, 114
 - location constraint, 112
 - meta-attribute, 114
 - node attribute, 114
 - node attribute expression, 109
 - operation defaults, 114
 - operation expression, 111

- options, 104
- resource defaults, 114
- resource expression, 111
- role, 112
- score, 112
- score-attribute, 112
- XML element, 104

run-command

- docker attribute, 96
- podman attribute, 96

S

SBD_SYNC_RESOURCE_STARTUP

- node option, 18

SBD_WATCHDOG_TIMEOUT

- node option, 18

score

- cluster_property_set, 20
- instance_attributes, 20
- meta_attributes, 20
- node health attribute value, 34
- resource_set attribute, 63
- rsc_colocation attribute, 60
- rsc_location attribute, 55
- rule, 112
- type, 11
- utilization, 20

score-attribute

- rule, 112

seconds

- date_spec, 106
- duration, 107

select

- XML element, 126

select_attributes

- XML element, 126

select_fencing

- XML element, 126

select_nodes

- XML element, 126

select_resources

- XML element, 126

sequential

- resource_set attribute, 63

server

- remote node, 44

shutdown

- node attribute, 33

shutdown-escalation

- cluster option, 29

shutdown-lock

- cluster option, 28

shutdown-lock-limit

- cluster option, 28

site-name

- node attribute, 33

source-dir

- storage-mapping attribute, 98

source-dir-root

- storage-mapping attribute, 98

standby

- node attribute, 33

start

- date_expression, 105

start-delay

- action property, 48
- op, 48

start-failure-is-fatal

- cluster option, 24

startup-fencing

- cluster option, 29

status

- XML element, 134

STONITH, 70

- resources, 38

stonith-action

- cluster option, 24

stonith-enabled

- cluster option, 24

stonith-max-attempts

- cluster option, 24

stonith-timeout

- cluster option, 24

stonith-timeout (primitive instance attribute), 72

stonith-watchdog-timeout

- cluster option, 25

stop-all-resources

- cluster option, 23

stop-orphan-actions

- cluster option, 24

stop-orphan-resources

- cluster option, 23

storage-mapping

- attribute, id, 98
- attribute, options, 98
- attribute, source-dir, 98
- attribute, source-dir-root, 98
- attribute, target-dir, 98

symmetric-cluster

- cluster option, 23

symmetrical

- rsc_order attribute, 59

Symmetrical Clusters, 57

System Service

- resources, 37

Systemd

- resources, 37

T

- target
 - fencing-level, 84
- target-attribute
 - fencing-level, 84
- target-dir
 - storage-mapping attribute, 98
- target-pattern
 - fencing-level, 84
- target-role
 - resource option, 39
- target-value
 - fencing-level, 84
- terminate
 - node attribute, 33
- text
 - type, 11
- then
 - rsc_order attribute, 58
- timeout
 - action property, 46
 - alert meta-attribute, 125
 - op, 46
 - type, 11
- timestamp-format
 - alert meta-attribute, 125
- transient_attributes
 - XML element, 135
- transition-delay
 - cluster option, 29
- transition-key
 - lrm_rsc_op, 136
- transition-magic
 - lrm_rsc_op, 136
- type
 - boolean, 10
 - date/time, 10
 - duration, 10
 - enumeration, 11
 - epoch_time, 11
 - expression, 110
 - id, 11
 - integer, 11
 - iso8601, 11
 - nonnegative integer, 11
 - percentage, 11
 - port, 11
 - range, 11
 - resource, 38
 - rsc_expression, 111
 - score, 11
 - text, 11
 - timeout, 11
 - version, 11

U

- uname
 - node_state, 134
- unfencing, 76
- utilization
 - id, 20
 - score, 20

V

- VALGRIND_OPTS
 - node option, 19
- validate-with
 - cib, 21
- value
 - expression, 110
- value-source
 - expression, 110
- version
 - type, 11

W

- weekdays
 - date_spec, 106
- weeks
 - date_spec, 107
 - duration, 107
- weekyears
 - date_spec, 107
- with-rsc
 - rsc_colocation attribute, 60
- with-rsc-role
 - clone ordering constraint, 60

X

- XML element
 - acl_group, 120
 - acl_permission, 119
 - acl_role, 118
 - acl_target, 119
 - acls, 118
 - alert, 123
 - alerts, 123
 - attribute, 126
 - bundle, 96
 - cib, 19
 - clone, 90
 - configuration, 10, 19
 - date_expression, 105
 - date_spec, 106
 - docker, 96
 - duration, 107
 - expression, 109
 - group, 88

- lrm, 135
- lrm_resource, 135
- lrm_resources, 135
- lrm_rsc_op, 136
- network, 97
- node_state, 134
- op_expression, 111
- podman, 96
- port-mapping, 97
- recipient, 124
- resource_set, 62
- role, 120
- rsc_colocation, 60
- rsc_expression, 111
- rsc_location, 55
- rsc_order, 58
- rule, 104
- select, 126
- select_attributes, 126
- select_fencing, 126
- select_nodes, 126
- select_resources, 126
- status, 134
- transient_attributes, 135
- xpath
 - acl_permission attribute, 119

Y

- yeardays
 - date_spec, 106
- years
 - date_spec, 107
 - duration, 107
- yellow
 - node health attribute value, 34